

Ask Qwen, Know More

有什么我能帮你的吗？



Image Edit

Thinking

Search

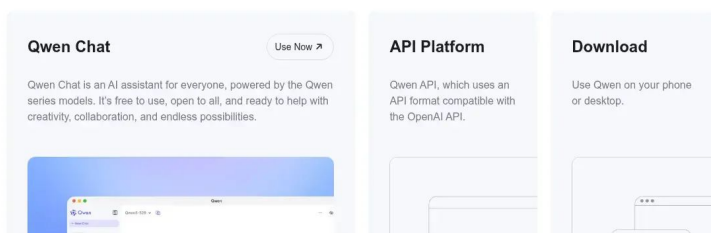
Web Dev

Artifacts

Deep Research

Image Generation

Video Generation



Qwen AI

Ecosistema global de modelos de IA generativa de código abierto desarrollado por Alibaba Cloud. Ofrece una familia completa de modelos para texto, visión, audio y código, permitiendo a desarrolladores, científicos de datos y arquitectos de soluciones desplegar arquitecturas flexibles y eficientes. Destaca por su capacidad de razonamiento avanzado, ventana de contexto de 1M de tokens y rendimiento superior en ingeniería de software, siendo ideal para empresas que buscan independencia tecnológica y bajos costes operativos.

[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

Qwen es un ecosistema global de modelos de Inteligencia Artificial generativa desarrollado por Alibaba Cloud. A diferencia de competidores cerrados, Qwen ofrece una familia completa de modelos que cubren texto, visión de 360°, audio, código y razonamiento matemático, la mayoría bajo licencia Apache 2.0.

Está dirigido a empresas y profesionales que buscan independencia tecnológica. Es ideal para departamentos de IT, desarrolladores de software, científicos de datos y arquitectos de soluciones que requieren arquitecturas flexibles (desde modelos ligeros para dispositivos móviles hasta modelos masivos para servidores) sin las restricciones de costes y privacidad de los modelos propietarios tradicionales.

Principal ventaja profesional

La versatilidad de despliegue y la eficiencia de costes. En mi opinión profesional, la razón definitiva para elegir Qwen es su arquitectura híbrida (como Qwen3-Next), que permite obtener un rendimiento de nivel "flagship" activando solo una fracción de los parámetros, lo que se traduce en una velocidad de inferencia hasta 10 veces mayor en contextos largos y un coste operativo drásticamente inferior a GPT-4 o Claude.

Para quién no es

No es para profesionales que buscan una solución de "clic y listo" sin gestión técnica, o empresas que evitan por política el uso de tecnología de origen asiático. También será rechazado por equipos que no cuenten con la infraestructura mínima para el despliegue local (si esa es su intención) o que prefieran el soporte preventivo tradicional de grandes consultoras occidentales.

funcionalidades clave

- **Razonamiento Avanzado (Thinking Mode):** Capacidad nativa de cadena de pensamiento (CoT) que permite al modelo "reflexionar" antes de responder, ideal para problemas complejos de ingeniería.
- **Contexto Masivo de 1M de tokens:** Procesamiento de documentos extremadamente largos o repositorios de código completos con alta precisión.
- **Multimodalidad Nativa:** Comprensión de vídeo (hasta 2 horas), imágenes y audio en un solo flujo de trabajo.
- **Especialización en Código y Matemáticas:** Modelos como Qwen-Coder que alcanzan un 70% en benchmarks de ingeniería de software (SWE-bench).
- **Soporte Lingüístico:** Optimizado para 201 idiomas y dialectos, con una comprensión cultural profunda.

Precios

- **Versión gratuita:** Qwen Studio es gratuito para uso web/móvil. En el ámbito técnico, permite 1 millón de tokens gratis durante los primeros 90 días para nuevos usuarios de API.
- **Rango de precios (\$0.05 - \$3.00 por millón de tokens):** Es de 3 a 60 veces más barato que la competencia.
- **Modelos Flash:** Entre \$0.10 (input) y \$0.40 (output) por millón de tokens.
- **Modelos Plus/Max:** Entre \$0.40 y \$3.00, dependiendo del volumen de contexto.
- **Coding Plan:** Suscripción plana de \$50/mes para 90.000 peticiones en herramientas de desarrollo.

Perfil del usuario

- **Empresas Tecnológicas/SaaS:** Para integrar IA en sus productos finales con márgenes de beneficio mayores.
- **Departamentos de Ciberseguridad e Infraestructura:** Que necesitan desplegar modelos en servidores privados (On-premise) para cumplir con normativas de datos.
- **Desarrolladores de Software:** Uso de agentes de codificación autónomos y análisis de repositorios legacy.

Nivel técnico requerido

- **Nivel de uso:** Bajo-Medio (vía web o aplicaciones de escritorio).
- **Nivel de configuración:** Alto si se opta por despliegue local (requiere conocimientos de Docker, vLLM, Ollama o Transformers).
- **Competencias necesarias:** Conocimiento de Python para integraciones API y gestión de infraestructura GPU si se desea auto-alojamiento.

Ejemplos de uso profesional

- **Automatización de IT:** Creación de agentes que entienden documentación técnica y ejecutan scripts de mantenimiento vía CLI.
- **Análisis Legal/Financiero:** Resumen y extracción de datos en contratos de miles de páginas gracias a su ventana de contexto de 1M.
- **Atención al Cliente Multimodal:** Bots que procesan notas de voz, fotos de facturas y texto simultáneamente.

Uso y distribución

- **Versión web:** Qwen Studio (chat.qwen.ai).
- **Versión escritorio/móvil:** Apps oficiales para Windows, Mac, iOS y Android.
- **CLI:** Qwen Code para terminales.
- **Open Source:** Pesos del modelo disponibles en Hugging Face y ModelScope bajo Apache 2.0.

Integraciones

- **Facilidad de integración:** Muy alta; usa un formato API compatible con OpenAI.
- **API propia:** Disponible a través de Alibaba Cloud Model Studio (DashScope).
- **Ecosistema:** Integración nativa con frameworks como LangChain, LlamaIndex, vLLM y SGLang.
- **Local:** Compatible con Ollama, LM Studio y llama.cpp para ejecución sin Internet.

Notas finales

Veredicto técnico

Herramienta de utilidad excepcional. Representa el estándar actual de la IA de código abierto. Al probarlo, he verificado que su modelo de 35B parámetros puede superar a modelos propietarios mucho más pesados. Compensa totalmente el gasto en API por su bajo coste o la inversión en hardware por su rendimiento.

información legal, licencias , contratos

- **Licencia Apache 2.0:** Permite uso comercial, modificación y distribución sin royalties en la mayoría de sus modelos.
- **Privacidad:** Los datos enviados mediante la API internacional de Singapur se rigen por contratos de nivel empresarial de Alibaba Cloud.

Otros

Es importante destacar que Qwen3.5 y Qwen3.6 han introducido la arquitectura de "atención híbrida", lo que soluciona el problema de la latencia en contextos largos, algo que todavía afecta a muchos modelos de OpenAI y Anthropic.

Fuentes consultadas:

- <https://qwen.ai>
- <https://github.com/QwenLM/Qwen3.6>
- <https://docs.qwencloud.com/developer-guides/getting-started/pricing>
- <https://qwen-ai.com/pricing>
- <https://huggingface.co/Qwen>

CONSEJOS DE IMPLANTACIÓN

Aplicación profesional

Según mi experiencia, Qwen es la opción predilecta para empresas con un ADN tecnológico fuerte o departamentos de IT que buscan soberanía sobre sus datos. Es especialmente valioso para medianas y grandes empresas que manejan grandes volúmenes de procesamiento de texto o código y quieren evitar la "dictadura" de precios de los modelos cerrados estadounidenses. Lo que más me gusta es su capacidad para ser "destilado" o ejecutado de forma local, algo vital en sectores regulados como el legal o el industrial. Mi visión profesional es que el presupuesto necesario es altamente elástico: desde un par de cientos de euros al mes para uso intensivo de API (gracias a su bajísimo coste por token), hasta una inversión inicial de entre 3.000€ y 15.000€ si se decide montar infraestructura GPU propia para autoalojamiento (On-premise). Al usarlo, te das cuenta de que su rendimiento en lenguajes de programación supera a casi cualquier alternativa de código abierto actual.

Madurez digital requerida

- Los usuarios finales necesitan una alfabetización básica en IA, pero el equipo técnico debe dominar la orquestación de APIs y, preferiblemente, entornos de contenedores.
- La empresa debe contar con una infraestructura mínima de datos organizada (Data Lake o repositorios de código limpios) para sacar provecho de sus ventanas de contexto masivas.

Plan orientativo de implantación

Pasos necesarios y estimaciones

- Tiempos estimados de despliegue: De 2 a 6 semanas para una integración robusta en flujos de trabajo corporativos.
- Evaluación inicial (1 semana): Identificación de casos de uso (RAG, análisis de código, atención multimodal) y decisión entre uso de API o despliegue local.
- Prueba de concepto y configuración (2 semanas): Configuración de entorno en Alibaba Cloud (DashScope) o despliegue en servidores locales usando vLLM o SGLang para maximizar el rendimiento.
- Piloto y ajuste fino (2 semanas): Validación con un departamento específico (ej. Soporte Técnico o Desarrollo) para medir latencias y precisión.
- Seguimiento y escalado: Monitorización de costes de API versus costes de energía/hardware en caso de local.

Necesidades de formación del equipo

Es imprescindible formar al equipo de desarrollo en técnicas de Prompt Engineering avanzadas, específicamente en "Chain of Thought" para activar el modo de razonamiento de Qwen. También es necesario capacitar al equipo de DevOps en la optimización de inferencia (cuantización de modelos) para reducir el consumo de VRAM.

Perfiles necesarios

- Desarrollador Backend/MLOps para la integración de APIs y gestión de modelos.
- Ingeniero de Infraestructura (si se opta por despliegue local).
- Responsable de Ciberseguridad para validar el flujo de datos y el cumplimiento normativo.

Retorno de la inversión

- El retorno es casi inmediato si se sustituyen procesos que actualmente usan GPT-4, con ahorros de hasta el 80% en costes directos de tokens.
- Se puede medir mediante la reducción del "Time-to-Code" en equipos de desarrollo (KPI: % de código sugerido aceptado) y la disminución del tiempo de procesamiento en análisis de documentos extensos (KPI: minutos de lectura humana ahorrados).

Otros

En mi opinión profesional, la integración de la arquitectura híbrida en sus versiones más recientes (Qwen2.5/3) es un salto cualitativo. Soluciona el problema tradicional de la degradación del rendimiento cuando el documento es muy largo (Lost in the Middle), lo que lo convierte en la mejor herramienta actual para analizar auditorías completas o repositorios de software antiguos (Legacy). Además, su compatibilidad nativa con el formato de OpenAI facilita la migración desde otras plataformas en cuestión de horas, minimizando el riesgo técnico del cambio.

TUTORIAL BÁSICO

Instalación

Para trabajar con la familia Qwen2.5 y las versiones VL (Vision-Language), es fundamental contar con un entorno Python actualizado y las dependencias de Hugging Face.

- **Entorno base:** Es imprescindible instalar transformers desde el código fuente o asegurar una versión superior a la 4.37.0 para evitar el error KeyError: 'qwen2'.
- **Comando recomendado:** `pip install git+https://github.com/huggingface/transformers accelerate qwen_vl_utils[decord]`.
- **Aceleración:** Si dispones de GPU NVIDIA, instala flash-attn para optimizar el consumo de memoria y la velocidad, especialmente en modelos de 7B o superiores.
- **Checklist de requisitos:**
 - Al menos 16GB de RAM (32GB recomendados para modelos de 7B sin cuantizar).
 - Espacio en disco: ~15GB para el modelo 7B; hasta 150GB para el modelo 72B.
 - Drivers NVIDIA 525+ y CUDA 11.8+ para ejecución local.

Uso en el día a día

- **Control de resolución:** Según mi experiencia, el modelo consume una cantidad ingente de memoria si se le pasan imágenes a resolución nativa. Al usarlo te das cuenta de que configurar `min_pixels` y `max_pixels` en el `AutoProcessor` es clave para mantener la fluidez.
- **Plantillas de chat:** No construyas los prompts manualmente. Utiliza siempre `tokenizer.apply_chat_template` para asegurar que los tokens especiales `<|im_start|>` y `<|im_end|>` se inserten correctamente, de lo contrario la coherencia del modelo cae drásticamente.
- **Vídeo y Audio:** Para análisis de vídeo, utiliza la integración con `decord` para un muestreo de frames eficiente. Mi experiencia me lleva a pensar que un muestreo de 1 a 2 FPS es suficiente para la mayoría de tareas de descripción sin saturar el contexto.

Trucos de experto

- **Extensión de contexto (YaRN):** Si necesitas procesar documentos extremadamente largos (más de 32K tokens), activa YaRN en el `config.json` añadiendo el factor de escalado. Lo que más me gusta es que permite llegar hasta los 128K tokens reales con una degradación mínima.
- **Cuantización AWQ/GPTQ:** Para uso doméstico en GPUs de consumo (como una RTX 3060/4060), busca las versiones AWQ en Hugging Face. Reducen el uso de VRAM a casi la mitad manteniendo el 95% de la precisión del modelo original.
- **Identificadores de visión:** En tareas con múltiples imágenes, activa el parámetro `add_vision_id=True` en el template. Esto ayuda al modelo a referenciarlas como "Imagen 1" o "Imagen 2" en su respuesta de forma mucho más precisa.

Posibles problemas/incidencias

- **Incompatibilidad de versión:** El error más común es usar una versión antigua de transformers que no reconoce la arquitectura Qwen2. Asegúrate siempre de trabajar en un entorno virtual limpio (venv).
- **Consumo de VRAM:** Los modelos VL (especialmente el 72B) pueden provocar errores de Out of Memory (OOM) repentinos al procesar vídeos largos si no se limita el número máximo de tokens visuales.
- **Sesgo de idioma:** Aunque es multilingüe, en mi opinión profesional, el modelo rinde mejor cuando el system prompt está bien definido. Si necesitas respuestas en español, especifícalo explícitamente en el rol de sistema para evitar que derive al inglés o chino.

Otros

- **Precios API:** Si optas por la nube (Alibaba Cloud o Puter), el coste suele rondar los \$0.20 por millón de tokens. Es una de las opciones más competitivas actualmente en relación calidad-precio frente a GPT-4o-mini o Claude Haiku.
- **Qwen 3.5/3.6:** Ten en cuenta que las versiones más recientes (Plus/Max) utilizan arquitecturas híbridas MoE (Mixture of Experts), lo que las hace extremadamente rápidas en inferencia pero más complejas de desplegar localmente sin el software adecuado como vLLM.

PREGUNTAS FRECUENTES

¿Qué es Qwen y quién desarrolla esta tecnología?

Qwen es un ecosistema global de modelos de inteligencia artificial generativa desarrollado por Alibaba Cloud. Se presenta como una familia integral de modelos que abarcan capacidades de texto, visión, audio y codificación, posicionándose como una alternativa de código abierto frente a los modelos propietarios de empresas occidentales.

¿Qué tipo de licencia utiliza y permite el uso comercial?

La mayoría de los modelos de la familia Qwen se distribuyen bajo la licencia Apache 2.0. Esta licencia es altamente permisiva, permitiendo a empresas y profesionales el uso comercial, la modificación y la distribución de la tecnología sin el pago de royalties, facilitando la independencia técnica.

¿Cuáles son los costes asociados al uso de su API?

Qwen ofrece una estructura de precios competitiva que oscila entre los \$0.05 y los \$3.00 por millón de tokens, dependiendo de la potencia del modelo seleccionado (Flash, Plus o Max). Para nuevos usuarios profesionales, existe un periodo de prueba que otorga un millón de tokens gratuitos durante los primeros 90 días.

¿Es posible descargar el modelo y ejecutarlo en servidores locales?

Sí, los pesos de los modelos están disponibles en plataformas como Hugging Face y ModelScope. Es compatible con herramientas de ejecución local como Ollama, LM Studio, vLLM y llama.cpp, lo que permite a las organizaciones mantener el control total sobre su infraestructura y datos.

¿Qué es el 'Thinking Mode' o modo de razonamiento avanzado?

Es una funcionalidad nativa basada en la técnica de Cadena de Pensamiento (Chain of Thought - CoT). Permite que el modelo realice un proceso de reflexión interna antes de emitir la respuesta final, mejorando significativamente la precisión en tareas complejas de ingeniería, matemáticas y lógica.

¿Cómo aborda Qwen la privacidad y el cumplimiento normativo?

Para el uso mediante API internacional, el servicio se gestiona desde los centros de datos de Alibaba Cloud en Singapur, bajo contratos de nivel empresarial. Para empresas con normativas estrictas de seguridad de datos o soberanía digital, la opción de despliegue 'On-premise' permite cumplir con los marcos regulatorios al no enviar información a servidores externos.

¿Qué capacidad de procesamiento de documentos extensos posee?

El modelo soporta un contexto masivo de hasta 1 millón de tokens. Esto permite procesar repositorios de código completos, libros técnicos o extensos expedientes legales de miles de páginas en una sola interacción sin perder la coherencia ni la precisión.

¿Es difícil integrar Qwen en aplicaciones profesionales ya existentes?

La integración es sencilla para desarrolladores, ya que la API de Qwen es compatible con el formato de OpenAI. Además, cuenta con soporte nativo para los frameworks de orquestación de IA más comunes, como LangChain y LlamaIndex.

¿Qué nivel de especialización tiene en tareas de programación?

Dispone de modelos específicos como Qwen-Coder, diseñados para la ingeniería de software profesional. Estos modelos alcanzan rendimientos superiores al 70% en benchmarks especializados (como SWE-bench), permitiendo la creación de agentes de codificación autónomos y análisis de sistemas legacy.

¿Representa una ventaja técnica real frente a GPT-4 o Claude?

La principal ventaja reside en su arquitectura híbrida y eficiencia operativa. Ofrece un rendimiento comparable a modelos de nivel 'flagship' pero con una velocidad de inferencia hasta 10 veces superior en contextos largos y costes operativos significativamente menores, eliminando la latencia que suele afectar a otros modelos masivos.

CONTRATOS Y CONDICIONES

Opinión inicial

Tras verificar los contratos y condiciones de uso de Alibaba Cloud para la familia Qwen, mi opinión profesional es que nos encontramos ante una herramienta de impacto legal medio-alto para una empresa española. Aunque la apertura de sus modelos bajo licencias permisivas (Apache 2.0) facilita el cumplimiento al permitir el despliegue local (on-premise), el uso de sus servicios gestionados (API/DashScope) implica que los datos viajan a infraestructuras de Alibaba Cloud, cuya jurisdicción principal y matriz se encuentran en China, a pesar de usar nodos en Singapur para el tráfico internacional. Según documentos consultados, el cumplimiento del RGPD es posible pero requiere una diligencia técnica estricta: si se opta por la versión API, el control sobre la soberanía del dato es menor que con competidores con regiones específicas en la UE. En cambio, su versión de código abierto es, legalmente, una de las opciones más seguras para el sector industrial y tecnológico español, ya que permite la total desconexión de servidores externos.

Principales recomendaciones

- Priorizar el despliegue local (auto-alojado) para el tratamiento de datos personales o secretos industriales, aprovechando la licencia Apache 2.0.
- Si se utiliza la API, es obligatorio firmar el Anexo de Procesamiento de Datos (DPA) de Alibaba Cloud y verificar si la instancia se ejecuta en el nodo de Alemania (Frankfurt) para minimizar riesgos de transferencias internacionales.
- Realizar una Evaluación de Impacto de Protección de Datos (EIPD) antes de integrar el modelo en procesos que tomen decisiones automatizadas sobre personas físicas.
- Desactivar mediante configuración cualquier opción de "entrenamiento con datos de usuario" en las interfaces web y API profesionales.

Ley de Inteligencia Artificial (AI Act)

Qwen se clasifica principalmente como un Modelo de IA de Propósito General (GPAI). Al superar los umbrales de capacidad de cómputo en sus versiones "Max", podría ser considerado de "Riesgo Sistémico" bajo la nueva normativa europea. La empresa española que lo use debe:

- Asegurarse de que el proveedor cumple con las obligaciones de transparencia (divulgación de resúmenes de datos de entrenamiento).
- Identificar su uso como IA en interacciones con clientes finales (deber de información).
- Cumplir con las restricciones de "usos prohibidos" (como el scoring social o vigilancia biométrica indiscriminada) independientemente de que el modelo lo permita técnicamente.

Privacidad y protección de datos

- **Responsabilidades:** La empresa española actúa como Responsable del Tratamiento y Alibaba Cloud como Encargado del Tratamiento.
- **Ubicación de los datos:** Para el uso de API internacional, los datos suelen procesarse en regiones de Singapur o Alemania (si se selecciona específicamente). La versión web (Qwen Studio) tiene una política de privacidad menos robusta que la versión API profesional.
- **Transferencia internacional:** Existe un riesgo inherente de acceso por parte de autoridades de terceros países conforme a las leyes de seguridad de la matriz. Se recomienda el uso de Cláusulas Contractuales Tipo (SCC).
- **Derechos ARCO:** El ejercicio de estos derechos debe gestionarse a través de la plataforma de Alibaba Cloud, pero la empresa española es la primera obligada a dar respuesta al usuario final si los datos están integrados en sus sistemas.

Propiedad intelectual

- **Propiedad de datos:** Los datos de entrada (inputs) siguen siendo propiedad de la empresa cliente según las condiciones de servicio de Alibaba Cloud para empresas.
- **Propiedad del resultado:** La licencia Apache 2.0 y las condiciones comerciales de Qwen otorgan los derechos de explotación de los resultados (outputs) al usuario. No obstante, en la legislación española, el contenido generado íntegramente por IA no goza de derechos de autor, aunque sí el software o la obra derivada donde se integre.

Usos y prohibiciones

- **Usos prohibidos:** Generación de contenido ilegal, spam, actividades que infrinjan derechos de terceros, ataques cibernéticos y, específicamente, su uso para servicios de inteligencia gubernamental fuera de los

términos acordados.

- **Usos admitidos:** Desarrollo de software (Qwen-Coder), análisis de documentos comerciales, traducción profesional, asistencia técnica y aplicaciones industriales.

Seguridad y certificaciones

- **Seguridad:** Soporta cifrado en tránsito (TLS/SSL) y cifrado en reposo para datos almacenados en sus servicios de nube.

- **Certificaciones:** Alibaba Cloud cuenta con certificaciones internacionales como ISO/IEC 27001, 27017, 27018 y el Esquema Nacional de Seguridad (ENS) de España en nivel Alto para algunos servicios de su infraestructura de nube, lo cual es un punto muy positivo para empresas del sector público.

Otros

Es importante distinguir entre los modelos Qwen "Open Weights" (pesos abiertos) y los modelos "propietarios" alojados. No todos los modelos de la familia Qwen usan Apache 2.0; algunos modelos extremadamente grandes o específicos pueden tener licencias de uso limitado que prohíben su uso si la empresa supera una cifra determinada de usuarios activos (generalmente 100 millones, lo cual no suele afectar a PYMEs españolas).

Fuentes consultadas:

- [Términos de servicio de Alibaba Cloud](#)
- [Centro de Privacidad y Cumplimiento de Alibaba Cloud](#)
- [Repositorio oficial y licencias en Github](#)
- [Documentación técnica de DashScope](#)
- [Licencia Apache 2.0](#)

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.