

# 第三方 API 模型接入清单

当前工作流程中通过 API 方式接入的全部第三方模型 —— 涵盖文本推理、图像生成、语音合成与数字人/视频四类。每一项标注模型名称、申请入口、能力价值与官方核实价格，便于横向对照与选型。

14 接入模型 / 服务    4 能力类别    9 厂商 / 平台    8/9 价格官方已核实

## 01 文本 LLM · 推理与对话

4 项

| 模型 / 厂商   | 特点 · 价值 · 场景  | 价格 · 申请  |
|---|---|--|
| <b>Claude (Anthropic)</b><br>Opus 4.8 · Sonnet 4.6<br>官方直连<br>● 主力          | 当前 <b>推理质量天花板</b> ，Agentic 编码与复杂多步任务最稳，工具调用可靠。Sonnet 享 1M 长上下文。<br>推理最强 Agentic 编码 长上下文 1M 工具调用稳<br>场景 Claude Code 核心引擎、复杂编码与多步推理 | <b>\$5</b> / \$25<br>Opus 4.8 输入 / 输出 (每百万 token)<br>Sonnet 4.6 <b>\$3 / \$15</b> ; 缓存读取 0.1×; Batch 5 折<br><a href="https://console.anthropic.com">console.anthropic.com</a>                    |
| <b>OpenRouter</b><br>LLM 聚合路由<br>300+ 模型统一网关<br>● 主力                        | 一个 Key <b>统一访问 300+ 跨厂商模型</b> ，自动故障转移与负载均衡，切换零成本，不被单一厂商锁定。<br>300+ 模型 一个 Key 自动容灾 不锁厂商<br>场景 多模型统一调度、对话与内容生成的中转层                  | <b>原价透传</b> 不加价<br>充值收 5.5% 手续费; BYOK 每月前 100 万次免费、之后 5%<br><a href="https://openrouter.ai">openrouter.ai</a>  |
| <b>Gemini Flash (Google)</b><br>2.5 / 3.x Flash<br>Google AI Studio<br>● 主力 | 高性价比多模态模型， <b>原生支持文件上传与视频分析</b> ，超大上下文。适合大批量、低延迟的理解类任务。<br>多模态 视频分析 超大上下文 便宜<br>场景 视频内容分析、大批量文本理解与分类                              | <b>\$0.30</b> / \$2.50<br>2.5 Flash 输入 / 输出 (每百万 token)<br>3.5 Flash <b>\$1.50 / \$9</b> ; Flash-Lite 低至 <b>\$0.10 / \$0.40</b><br><a href="https://aistudio.google.com">aistudio.google.com</a> |
| <b>Qwen3 (SiliconFlow)</b><br>Qwen3.5 系列<br>硅基流动 · 国内直连<br>● 主力             | 国内直连的 <b>开源 Qwen3 托管推理</b> ，单价极低、低延迟，适合高频小任务。新账号含免费额度。<br>国内直连 单价极低 低延迟 开源模型<br>场景 意图分类、关键词抽取等高频轻量任务                              | <b>¥0.40</b> 起 / 百万 token<br>输入 ¥0.40-1.20 (0-128k 档)<br>输出 <b>¥3.20-7.20</b> ; 长上下文档约 2.5-4×<br><a href="https://siliconflow.cn">siliconflow.cn</a>   |

## 02 图像生成 · 文生图与编辑

6 项

| 模型 / 厂商   | 特点 · 价值 · 场景   | 价格 · 申请   |
|---|--|---|
| <b>Gemini 3 Pro Image</b><br>gemini-3-pro-image-preview<br>Google · nano-banana-pro<br>● 主力       | <b>原生 4K + 最高保真度</b> ，多参考图严格一致性，中文文字渲染最佳。40s 出 4K 不需 upscale。<br>原生 4K 多参考图一致 中文文字准 色彩还原强<br>场景 商业级成片、4K 电商产品图、隐形模特  | <b>≈\$0.13</b> / 张 (1-2K)<br>4K 约 <b>\$0.24</b> /张; 输入 \$2、图像输出 \$120 (每百万 token); Batch 半价<br><a href="https://aistudio.google.com">aistudio.google.com</a>                      |
| <b>Gemini 3.1 Flash Image</b><br>gemini-3.1-flash-image-preview<br>Google · nano-banana 2<br>● 主力 | 1024 原生分辨率， <b>快且便宜</b> ，中文渲染好。Pro 版的轻量替身，用于不追 4K 的高频出图。<br>快 便宜 中文好 1024 原生<br>场景 中文 slogan 配图、内容快稿迭代               | <b>≈\$0.045</b> / 张 起<br>\$0.045 (0.5K) → <b>\$0.151</b> (4K) 按分辨率<br>输入 \$0.50、图像输出 \$60 (每百万 token)<br><a href="https://aistudio.google.com">aistudio.google.com</a>            |
| <b>gpt-image-2 (OpenAI)</b><br>gpt-image-2<br>OpenAI 官方直连<br>● 主力                                 | <b>指令跟随与图像编辑最强</b> ：保留原图主体换场景、重建光影逻辑。edits 最多 16 张参考图。<br>编辑最强 产品一致性 光影逻辑 16 张参考图<br>场景 产品换场景、风格迁移、角色跨镜一致性           | <b>≈\$0.05</b> / 张 (1024 <sup>2</sup> )<br>低 \$0.006 · 高 <b>\$0.21</b> ; 图像输入 \$8、输出 \$30 (每百万 token); Batch 5 折<br><a href="https://platform.openai.com">platform.openai.com</a> |
| <b>即梦 Seedream 4.0</b><br>doubao-seedream-4-0<br>火山方舟 Ark · 国内直连<br>● 可用                          | <b>中文理解与文字渲染强</b> ，多图融合 (最多 10 张)、序列图 (最多 15 张)。国内直连、单张成本最低。<br>中文准 多图融合 10 序列图 15 国内最便宜<br>场景 中文商业素材、批量风格统一配图       | <b>¥0.20</b> / 张<br>按输出图片计费 (RMB)<br>商业素材务必走 API，不用 Web UI<br><a href="https://console.volcengine.com/ark">console.volcengine.com/ark</a>   |
| <b>fal.ai</b><br>FLUX / Seedream / Nano Banana<br>多模型聚合平台<br>● 备用 · 余额耗尽                          | 一个 API <b>聚合多家开源/托管图像视频模型</b> (含 Recraft 矢量化)，按产出计费，排队与服务端错误不收费。<br>多模型聚合 矢量化 按产出计费 错误不收费<br>场景 跨模型快速试验、矢量 logo、备用通道 | <b>\$0.03</b> / 张 起<br>Seedream V4 \$0.03 · Nano Banana \$0.0398 · FLUX Kontext Pro <b>\$0.04</b> /张<br><a href="https://fal.ai/pricing">fal.ai/pricing</a>                       |
| <b>可灵 Kling</b><br>Kling · Kolors 可图<br>快手可灵开放平台<br>● 暂停 · 无余额                                    | <b>视频生成业界领先</b> (运动幅度 / 时长)，图像生成 (可图 Kolors) 中英文文字渲染强，面向 B2B 集成。<br>视频运动强 长时长 中英文文字准 B2B<br>场景 高动态视频生成 (图像为辅)        | <b>\$1</b> = 66 积分<br>≈\$0.015/积分; 视频 9-16 积分/秒<br>图像单价 <b>未公开核实</b> (站点拦截抓取)<br><a href="https://klingai.com/global/dev">klingai.com/global/dev</a>                              |

## 03 语音 TTS · 合成与克隆

2 项

| 模型 / 厂商  | 特点 · 价值 · 场景   | 价格 · 申请  |
|--|--|--|
| <b>MiniMax TTS</b><br>speech-2.8-hd<br>MiniMax · 国内直连<br>● 主力                                | <b>高拟真、情感丰富的中文 TTS</b> ，当前最高质量档。支持音色克隆、timber 混音 (中文音色说地道英文) 与跨语种。<br>最高拟真 音色克隆 混音 跨语种<br>场景 微信语音回复、视频号口播、数字人配音              | <b>¥3.5</b> / 万字符<br>turbo 档 <b>¥2</b> /万字符; 按输入字符计, 1 汉字 = 2 字符<br><a href="https://platform.minimaxi.com">platform.minimaxi.com</a>            |
| <b>Qwen3-TTS (阿里云百炼)</b><br>qwen3-tts-flash · cosyvoice-v3-plus<br>DashScope · 国内直连<br>● 迁移中 | <b>17 种方言 + 10 种外语</b> ，音色克隆 1 年有效 (MiniMax 仅 7 天)，自然语言情绪控制。单价低于 MiniMax。<br>17 方言 克隆 1 年 情绪控制 更便宜<br>场景 微信语音 (已上线)、视频号口播替换中 | <b>¥0.115</b> / 万字符<br>qwen3-tts-flash 国内价; 音色复刻 <b>\$0.01</b> /个<br><a href="https://bailian.console.aliyun.com">bailian.console.aliyun.com</a> |

## 04 数字人 / 视频 · 生成

2 项

| 模型 / 厂商  | 特点 · 价值 · 场景  | 价格 · 申请   |
|--|---|---|
| <b>石榴数字人</b><br>16ai · createByAudioFile<br>向量方程科技<br>● 可用                               | <b>一站式数字人</b> : 声音合成 + 形象训练 (照片/视频建模) + 音频驱动视频生成。覆盖广告、医疗、教育等行业。<br>一站式 音频驱动 照片建模 多行业<br>场景 数字人口播视频 (TTS 出音 → 驱动形象)          | <b>¥6</b> / 分钟<br>数字人视频生成 <b>¥6</b> / 分钟<br><a href="https://api.16ai.vip">api.16ai.vip</a>   |
| <b>即梦视频 Seedance</b><br>doubao-seedance-1.0-pro / fast / lite<br>火山方舟 Ark · 国内直连<br>● 可用 | <b>文生 / 图生视频</b> 。1.0-pro 高画质，pro-fast 提速 3 倍降价 72% 量产，lite 入门档。国内直连。<br>文生/图生视频 高画质 pro 量产 fast 国内直连<br>场景 短剧、水墨漫剧、视频号批量出片 | <b>¥3.67</b> / 5s 1080p<br>1.0-pro ¥15/百万 token; pro-fast <b>≈¥1.03/5s</b> ;<br>离线 batch 5 折<br><a href="https://console.volcengine.com/ark">console.volcengine.com/ark</a> |

● 主力 / 可用 —— 当前在用    ● 迁移中 —— 逐步替换    ● 备用 / 暂停 —— 余额或额度问题

价格说明 · 上表价格于 2026-06-04 核实自各厂商官方定价页，单位与币种以原页面为准 (USD / RMB 混用，已标注)。除可灵图像单价 (站点拦截抓取) 外，其余均锚定官方页面。

价格随官方政策动态调整，实际计费以下单时官网为准。缓存、Batch、长上下文等折扣档详见各官方文档。