

MiniMax-M1: 智能计算的新突破

如何高效处理复杂任务并实现卓越性能



什么是 MiniMax-M1 ?

🌐 全球首个开源混合注意力模型

MiniMax-M1是世界上第一个开源的大规模混合注意力推理模型，为AI领域带来突破性进展。

⚙️ 混合专家模型架构

采用混合专家模型(MoE)架构，总计4560亿参数，每个token激活459亿参数，配合32个专家系统。

⚡ 闪电注意力机制

创新的闪电注意力机制大幅提升计算效率，使模型能够处理超长文本。

🔲 超长上下文支持

原生支持100万tokens的上下文长度，是DeepSeek R1的8倍，远超其他开源模型。



MiniMax-M1为何如此高效？

⚠️ 传统模型的挑战

传统Transformer架构在处理长文本时，注意力机制的计算复杂度呈**二次方增长**，导致计算成本急剧上升。

💡 闪电注意力解决方案

MiniMax-M1采用创新的闪电注意力机制，每7个闪电注意力块后接一个传统注意力块，实现**线性计算复杂度**。

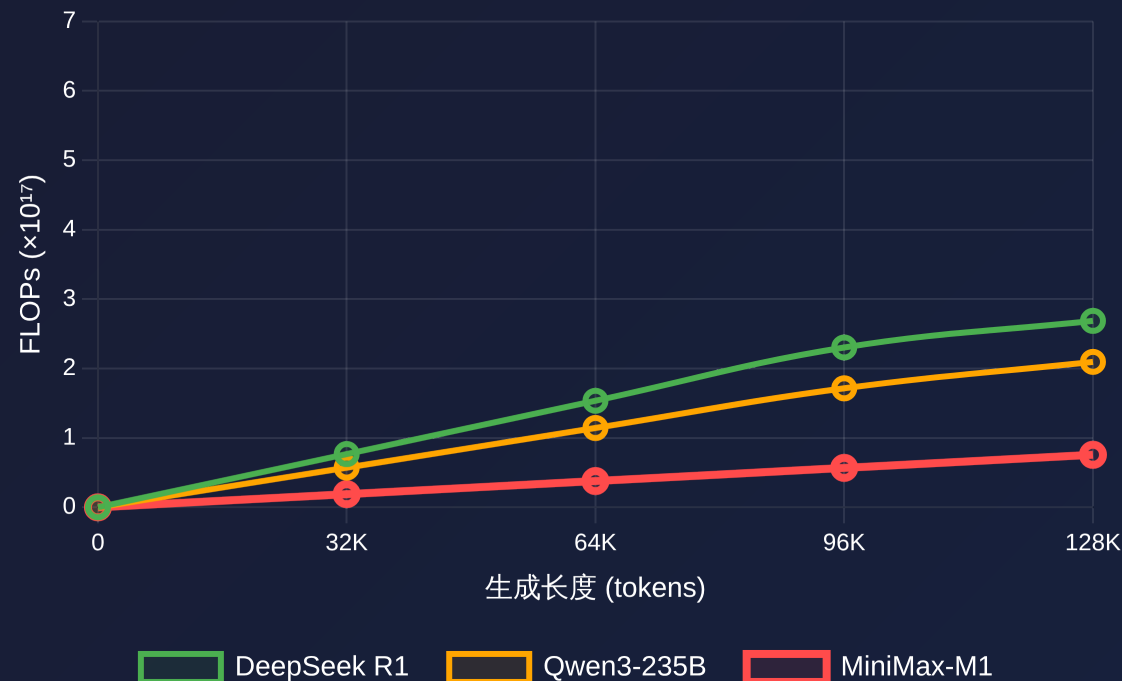
📈 显著的效率提升

在10万tokens生成长度下，MiniMax-M1仅消耗DeepSeek R1的**25%计算量**，在6.4万tokens时消耗不到50%。

⚡ 适用于复杂任务

高效的计算特性使MiniMax-M1特别适合处理**需要长输入和大量思考**的复杂任务，如数学推理和软件工程。

计算量随生成长度的增长对比



MiniMax-M1的计算效率优势随着生成长度增加而更加明显

MiniMax-M1是如何训练的？

🧠 大规模强化学习

MiniMax-M1通过**大规模强化学习(RL)**训练，使模型能够进行长时间推理，解决复杂问题。

🔗 创新CISPO算法

提出**CISPO算法**，通过裁剪重要性采样权重而非token更新，比其他RL变体效率更高，实现2倍速度提升。

📊 多样化训练数据

训练数据涵盖**数学推理**、**竞争性编程**、**逻辑推理**等可验证问题，以及问答和创意写作等不可验证问题。

训练流程

- 1 持续预训练：在7.5T高质量token上继续训练MiniMax-Text-01模型
- 2 监督微调：注入思维链(CoT)模式，为强化学习奠定基础
- 3 强化学习：使用CISPO算法进行高效训练

512

H800 GPU

3周

训练时间

\$53.47万

训练成本

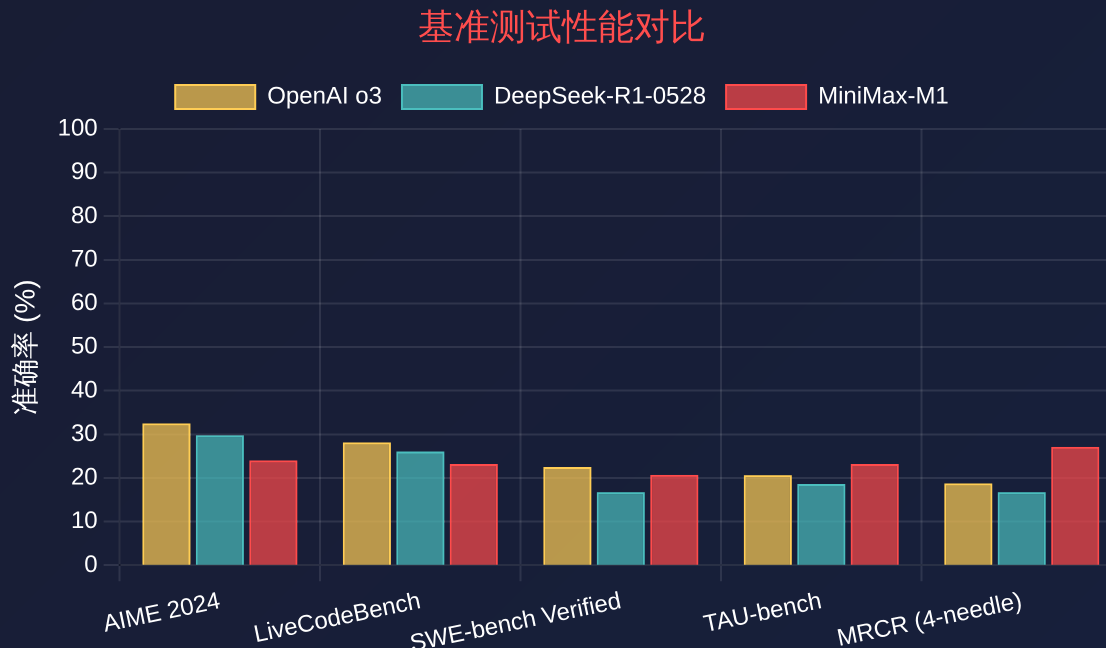


MiniMax-M1的性能表现

MiniMax-M1在多项基准测试中展现出**卓越性能**，尤其在复杂软件工程、工具使用和长上下文任务方面表现突出。

与其他开源模型相比，MiniMax-M1在整体上**超越了原始DeepSeek-R1和Qwen3-235B**，在某些领域甚至可与闭源商业模型媲美。

特别值得注意的是，MiniMax-M1在**TAU-Bench**（代理工具使用基准）上超越了Gemini 2.5 Pro，在**长上下文理解**基准上超越了OpenAI o3和Claude 4 Opus。



数据来源：MiniMax-M1-80k模型在各项基准测试中的表现



软件工程能力

在SWE-bench Verified测试中表现优异，能够解决复杂的软件工程问题和bug修复任务。



代理工具使用

在TAU-bench测试中超越Gemini 2.5 Pro，展示出卓越的工具使用和任务自动化能力。



长上下文理解

在MRCR (4-needle)测试中表现突出，能够有效处理和理解超长文本内容。

总结与展望

MiniMax-M1通过创新的**混合注意力架构**和**闪电注意力机制**，实现了测试时计算的高效扩展，为下一代语言模型代理奠定了坚实基础。

其**100万tokens的超长上下文**支持和**显著降低的计算成本**，使其特别适合解决需要处理长输入和大量思考的复杂现实世界挑战。

- ✓ MiniMax-M1在**复杂软件工程、工具使用和长上下文任务**方面表现尤为突出，展示了其作为通用AI助手的潜力。
- ✓ 通过**CISPO算法**的创新，MiniMax-M1实现了更高效的强化学习训练，大幅降低了训练成本和时间。
- ✓ MiniMax-M1的**开源发布**将促进AI领域的合作与进步，推动大型推理模型的广泛应用。



获取与使用

-  GitHub开源代码库
-  支持vLLM和Transformers框架
-  提供商业标准API
-  详细部署指南与文档

访问 minimax.io 了解更多信息