

【嬴政指数】2026 Q1 硬件天梯图

# DeepSeek 最佳运行硬件排行榜

深度评测报告：从数据中心到你的口袋，谁才是穷人的法拉利？

实测 RTX 4090 vs Mac M4 vs 小米8，AI 推理性能全面横评

NVIDIA RTX 4090

NVIDIA A100 80GB

Apple Mac M4 Max

小米8 (骁龙845)

测试指标：Token/s 生成速度 | 显存占用 | 发热 & 功耗 | 综合性价比

独立测试机构：Winzheng Research Lab | 嬴政研究院

发布时间：2026年2月 | 基于社区实测数据与专业基准测试综合整理

# 一、开篇：为什么这份排行榜很重要？

2025年初，DeepSeek-R1 的发布彻底改变了 AI 行业的格局。这个拥有 671 亿参数的开源推理模型，在多项基准测试中达到了与 OpenAI o1 匹敌的水平，却以开源、免费的方式提供给全球开发者。这意味着一件事：你不再需要付费订阅才能用到顶级 AI，但你需要合适的硬件来跑它。

但问题来了：到底什么硬件能跑 DeepSeek？一张 RTX 4090 够不够？Mac 用户有没有机会？手里握着一部旧手机的穷学生，还有没有希望体验 AI 的魅力？

Winzheng Research Lab 将用实测数据告诉你答案。我们跨越了从数据中心级硬件到你口袋里的旧手机，测试了四款具有代表性的设备，并从生成速度、显存占用、发热/功耗、性价比等多个维度给出终极排名。

## 二、测试方法论

### 2.1 测试对象一览

设备	类型	显存/内存	带宽	参考价格	定位
NVIDIA A100 80GB	数据中心GPU	80GB HBM2e	2.0 TB/s	~\$2/hr (云)	专业旗舰
NVIDIA RTX 4090	消费级GPU	24GB GDDR6X	1.0 TB/s	¥~13,000	极客之王
Mac M4 Max 128GB	Apple Silicon	128GB 统一内存	546 GB/s	¥~28,000	创意专业人士
小米8 (骁龙845)	旧款手机	6GB RAM	~30 GB/s	¥~300 (二手)	极限挑战

### 2.2 测试模型与工具

由于完整版 DeepSeek-R1 671B 需要超过 400GB 显存，单卡无法直接运行，因此我们根据硬件能力选择了不同规格的蒸馏版模型进行测试：

A100 80GB : DeepSeek-R1-Distill-Qwen-32B (FP16) / 671B 1.58-bit 量化版

RTX 4090 : DeepSeek-R1-Distill-14B (Q4) / Distill-32B (Q4)

Mac M4 Max 128GB : DeepSeek-R1-Distill-70B (Q4) / Distill-32B (Q4)

小米8 : DeepSeek-R1-Distill-1.5B (Q4) — 唯一可运行版本

测试工具：Ollama + llama.cpp，Mac 额外使用 MLX 框架，手机端使用 Termux + Ollama。所有测试均在室温 25°C 环境下多次运行取平均值。

## 三、核心测试结果

### 3.1 同模型横向对比 (R1-Distill-14B Q4)

R1-Distill-14B Q4	A100 80GB	RTX 4090	Mac M4 Max	小米8
生成速度 (Token/s)	~85	~58	~22	N/A
显存/内存占用	~9GB / 80GB	~9GB / 24GB	~9GB / 128GB	无法加载
满载功耗	~250W	~350W	~45W	—
温度表现	风冷服务器	GPU ~72°C	微温 ~42°C	—

### 3.2 各设备最佳工作区间

设备	最佳模型	Token/s	显存占用	功耗	体验评级
A100	Distill-32B FP16	~45	~64GB	~300W	★★★★★
4090	Distill-14B Q4	~58	~9GB	~350W	★★★★★
4090	Distill-32B Q4	~34	~20GB	~350W	★★★★
M4 Max	Distill-70B Q4	~10	~40GB	~45W	★★★
M4 Max	Distill-32B Q4	~14	~20GB	~40W	★★★★
小米8	Distill-1.5B Q4	~3-5	~1.5GB	~5W	★☆

## 四、逐机深度解析

### 4.1 NVIDIA A100 80GB — "数据中心的霸主"

A100 是 NVIDIA Ampere 架构旗舰，80GB HBM2e 显存配合高达 2.0 TB/s 的内存带宽，是数据中心 AI 推理的主力。在 Dual A100 配置下运行 DeepSeek-R1 70B 模型，实测达到约 19.34 Token/s 的评估速率，GPU 利用率稳定在 ~88%。而运行蒸馏版 7B/14B 模型时，单卡吐出轻松超过 80 Token/s。

核心优势：巨大的显存容量和极高带宽，可以  
参数模型，无需量化压缩，保持全部模型精度。在批处理场景下表现尤其出色。

成本考量：云服务器租用价格约  
\$10,000-\$15,000。对个人用户不现实，但对企业用户是性能天花板。

结论：A100 是「不差钱」玩家的终极武器。但对普通用户而言，这是一辆「照字面意义的法拉利」——  
好，但太贵。

## 4.2 NVIDIA RTX 4090 — "极客之王"

RTX 4090 是消费级 GPU 的绝对王者。24GB GDDR6X 显存配合约 1.0 TB/s 带宽，让它在本地推理场景中表现出色。运行 14B 参数模型时达到约 58 Token/s，这比大多数人的阅读速度还快。即便是 32B 模型也能达到 ~34 Token/s，相当于 H100 约 75% 的性能，成本仅为后者的 21%。

显存限制：24GB 显存是最大短板。只能跑 Q4 量化的 32B 模型 (~20GB 显存占用)，70B 模型必须借助系统内存做混合推理，速度会大打折扣。对于完整版 671B，即使 1.58-bit 极限量化，单卡也只能卸载 ~7 层到 GPU。

发热功耗：满载推理约 350W，GPU 温度 ~72°C，风扇噪音明显但可接受。需要良好的机箱风道。

**结论：RTX 4090**  
是我们心目中的「穷人的法拉利」最强候选人。花万元出头，就能在本地获得媲美云服务的体验。

## 4.3 Apple Mac M4 Max 128GB — "沉默的极客"

Apple Silicon 的统一内存架构是它在 AI 推理领域的杀手锏。M4 Max 的 128GB 统一内存让它能运行完整的 70B 参数 Q4 量化模型，这是单张 RTX 4090 根本做不到的。70B 模型约 10 Token/s，32B 模型约 14 Token/s。

能效之王：整机功耗仅 ~40-45W，每瓦特产出的 Token 数远超 GPU 方案。不需要独显、不需要大功率电源、运行时几乎无声。对于长时间运行的开发场景，这是巨大优势。

生态系统：MLX 框架为 Apple Silicon 提供专门优化，性能优于通用 Ollama/llama.cpp 方案。Apple 生态的成熟度和易用性也是加分项。

**结论：如果你已是 Mac 用户且内存足够，M4 Max 是「沉默的极客」——看起来像普通笔记本，实际能跑 70B 大模型。**

## 4.4 小米8 (骁龙845) — "极限生存挑战"

小米8 发布于 2018 年，骁龙 845 处理器，6GB RAM，二手约 300 元。通过 Termux + Ollama 成功加载 DeepSeek-R1-Distill-1.5B，生成速度约 3-5 Token/s，勉强处于可阅读边缘。但要注意，这是 1.5B 参数的极小模型，推理能力和完整版 R1 有天壤之别。

发热情况：持续推理时机身温度迅速攻上 45°C+，局部热点接近 50°C。CPU 开始降频，生成速度从 5 降至 2-3 Token/s。手机变成「暖手宝」。

**警告：**请勿在充电时或覆盖物下进行测试，除非你想亲身体验「AI 发烧友」的字面含义。

结论：小米8 证明了「旧手机也能跑 AI」，但这更像是一种行为艺术（Performance Art）——证明了可能性，但不适合日常使用。

## 五、终极排行榜

### 5.1 综合性能排名 (The Winzheng Tier List)

排名	设备	速度	容量	能效	易用	性价比	总分
1	RTX 4090	9.5	6.0	5.0	8.0	9.5	38.0
2	M4 Max	6.0	9.0	10.0	9.5	7.0	41.5
3	A100 80GB	10.0	10.0	4.0	5.0	3.0	32.0
4	小米8	1.5	1.0	7.0	3.0	2.0	14.5

注：各项满分10分。M4 Max 虽然总分高，但主要胜在容量和能效；若只看纯推理速度和性价比，4090 依然是当之无愧的王者。

### 5.2 谁才是「穷人的法拉利」？

#### 穷人的法拉利：NVIDIA RTX 4090

综合性价比第一 | 58 Token/s @ 14B | 万元级硬件 ≈ 云服务体验

为什么是 RTX 4090？因为它完美地平衡了性能、价格和可获得性。在本地运行 DeepSeek 蒸馏版时，它能给你接近云服务的体验，却不需要持续付费，数据完全私有。花一万多元，你就能拥有一台真正属于自己的 AI 工作站。

最佳能效比 Mac M4 Max 每瓦特Token产出 ≈ 4090的7倍	最强绝对性能 A100 80GB 可加载最大模型，速度无上限	最佳勇气奖 小米8 300元手机也能跑AI，虽然很烫
---	--------------------------------------	----------------------------------

## 六、购买建议：不同预算怎么选？

预算	推荐硬件	可跑模型	适合人群
300元以下	二手旧手机 6GB+	1.5B-3B	好奇心驱动的学生，尝鲜体验
3K-6K元	二手 RTX 3090	7B-14B Q4	AI 爱好者，本地跑中等模型
1-1.5万	RTX 4090 24GB	14B-32B Q4	AI 开发者，追求性价比极致
2-3万	Mac M4 Max 128GB	32B-70B Q4	Mac 用户，重视静音和大模型加载能力
10万+	Mac Studio M3 Ultra 512GB	671B Q4 完整版	企业用户，本地部署完整模型

## 七、结语：AI 民主化的时代已经到来

DeepSeek-R1

的开源发布，让顶级

AI

推理能力第一次真正走入了普通人的硬件。从我们的测试中可以看到：

1. 一张 RTX 4090 就能让你在本地获得接近云服务的 AI 体验；
2. Mac 用户凭借统一内存架构，可以运行远超同价位 GPU 的大模型；
3. 即便是 7 年前的旧手机，也能跑小规模 AI 模型；
4. AI 硬件的门槛正以前所未有的速度降低。

我们正在进入一个 AI 民主化的时代。不论你的预算是 300 元还是 3 万元，都能找到属于自己的「法拉利」。区别只在于，有的是真正的 F40，有的是比例模型——但它们都能让你感受到速度与激情。

### 免责声明

本报告数据基于 Winzheng Research Lab 实测、专业评测机构报告及开源项目基准测试综合整理。实际性能可能因软件版本、量化方法、系统配置等因素有所差异。价格信息以 2026 年 2 月市场价为准。

数据来源：NVIDIA 官方博客、DatabaseMart GPU Benchmark、llama.cpp 社区、Unsloth 实测数据、MacRumors 社区测试、Android Authority 手机 LLM 测试等。

解码智能，定义价值。

Focus: AI Benchmarks, Model Security, Hardware Optimization

Stance: 100% Independent & Objective.