

Mac Studio for Local AI: Is It Worth the Price?

February 26, 2026 · by Mark Bartlett

[Download this guide as PDF](#)

Quick Answer: The Mac Studio M4 Max with 128GB (\$3,699) runs 70B models at 30-45 tok/s via MLX and fits on a desk. The M3 Ultra with 256GB (\$5,999+) loads models that no single consumer GPU can touch, including Qwen3 235B. For tok/s per dollar, a dual RTX 3090 PC wins. For memory per dollar, silence, and always-on serving, the Mac Studio wins. Buy the M4 Max 128GB if you need 70B models and quiet operation. Buy the M3 Ultra if you need 70B at high quantization or want to run 100B+ models. Don't buy either if 14B is enough for your work – a MacBook Pro with 24GB does that for half the price.

 **More on this topic:** [Best Local LLMs for Mac 2026](#) · [Running LLMs on Mac M-Series](#) · [Ollama on Mac: Setup & Optimization](#) · [VRAM Requirements](#)

The Mac Studio is Apple's answer to a question most PC builders never ask: what if you could run a 70B language model from something the size of a thick paperback, with no fan noise, pulling 20 watts at idle?

It's not cheap. The AI-relevant configurations start around \$2,800 and go past \$10,000. An equivalent PC build with used RTX 3090s generates tokens faster for less money. So why would anyone buy a Mac Studio for AI?

Because memory. The M4 Max goes up to 128GB of unified memory. The M3 Ultra goes up to 512GB. No consumer GPU on earth offers that. A 70B model at Q8 quantization, a 100B+ MoE model, multiple models loaded simultaneously – these are things that require multi-GPU setups costing thousands on PC, if they're possible at all. On Mac Studio, you plug it in and pull the model.

Current Mac Studio lineup (2025)

Two chip families, two very different price brackets:

Config	Chip	CPU/GPU cores	Memory	Bandwidth	Starting price
Base M4 Max	M4 Max	14-core / 32-core	36 GB	410 GB/s	\$1,999
Upgraded M4 Max	M4 Max	16-core / 40-core	64-128 GB	546 GB/s	~\$2,799+

Config	Chip	CPU/GPU cores	Memory	Bandwidth	Starting price
Base M3 Ultra	M3 Ultra	28-core / 60-core	96 GB	819 GB/s	\$3,999
Upgraded M3 Ultra	M3 Ultra	32-core / 80-core	256-512 GB	819 GB/s	~\$5,999+

The M3 Ultra is essentially two M3 Max chips fused together via UltraFusion. Double the cores, double the memory bandwidth, double the maximum memory. The tradeoff is price and power draw (215W max vs ~120W for the M4 Max).

Note: there's no M4 Ultra yet. If you need more than 128GB, the M3 Ultra is your only Mac Studio option right now.

Which chip to pick

Skip the base M4 Max (36GB). It handles 14B models fine, but a \$1,199 Mac Mini M4 Pro with 48GB does the same job. You're buying a Mac Studio for the memory headroom, so start at 64GB minimum.

The decision tree:

- **M4 Max 64GB (\$2,799):** Runs 32B models comfortably. Good entry point if you're testing larger models but don't need 70B daily.
- **M4 Max 128GB (\$3,699):** The sweet spot for most AI work. Runs 70B models fully in memory at Q4-Q5. Fast enough for interactive use.
- **M3 Ultra 96GB (\$3,999):** Oddly positioned. More bandwidth than the M4 Max 128GB but less memory. Only makes sense if you value throughput on 32B models over the ability to load 70B.
- **M3 Ultra 256GB (\$5,999+):** 70B models at Q8 quality, multi-model setups, or 100B+ MoE models. This is where Mac Studio has no PC equivalent.
- **M3 Ultra 512GB (\$8,000+):** Research-grade. Multiple 70B models loaded at once, or the largest open models like Qwen3 235B-A22B at usable quantization.

What you can actually run

Raw specs mean nothing without benchmarks. Here's what each configuration handles in practice:

M4 Max 128GB

Model	Quant	Memory used	Speed (MLX)	Speed (Ollama)
Qwen 3 8B	Q4_K_M	~5 GB	~58 tok/s	~45 tok/s
Qwen 3 14B	Q4_K_M	~9 GB	~40 tok/s	~32 tok/s
Qwen 3 32B	Q4_K_M	~19 GB	~28 tok/s	~22 tok/s
Llama 3.3 70B	Q4_K_M	~40 GB	~15 tok/s	~10 tok/s
Llama 3.3 70B	Q8_0	~72 GB	~9 tok/s	~7 tok/s
Qwen3 235B-A22B	Q2_K	~88 GB	~10 tok/s	~8 tok/s

The 128GB config leaves ~40-50GB free after loading a 70B Q4 model. That's enough for macOS, a browser, and your development tools without hitting memory pressure.

The M4 Max's 546 GB/s bandwidth is noticeably faster than the base M4 Max (410 GB/s). Always get the 16-core/40-core variant if you're going 64GB+.

M3 Ultra 256GB

Model	Quant	Memory used	Speed (MLX)
Llama 3.3 70B	Q4_K_M	~40 GB	~20 tok/s
Llama 3.3 70B	Q8_0	~72 GB	~12 tok/s
Llama 3.3 70B + Qwen 3 32B	Both Q4	~59 GB	Both loaded, switch instantly
Qwen3 235B-A22B	Q4_K_M	~140 GB	~8 tok/s
DeepSeek V3 (MoE)	Q2	~180 GB	~5 tok/s

The M3 Ultra's 819 GB/s bandwidth pushes tokens about 50% faster than the M4 Max at the same model size. Where it really shines is headroom: loading two or three models simultaneously for a router/specialist setup, or running the largest available open models.

What only Mac Studio can do

A few things no single-GPU PC can match:

- **70B at Q8 quality** needs ~72GB. No consumer GPU has that. The M4 Max 128GB handles it.
- **Multi-model serving** — load a fast 8B router + a 32B specialist + a coding model. The 256GB M3 Ultra handles all three without swapping.

- **100B+ MoE models** — Qwen3 235B-A22B (22B active parameters, 235B total) runs at 8-10 tok/s. Needs ~88-140GB depending on quantization.
- **Fine-tuning with MLX** — unified memory means no VRAM wall. You can LoRA fine-tune a 14B model on the M4 Max 128GB without the out-of-memory crashes that plague 24GB GPUs.

Cost comparison: Mac Studio vs PC builds

This is where the conversation gets honest. Dollar for dollar, NVIDIA generates more tokens per second.

M4 Max 128GB (\$3,699) vs dual RTX 3090 PC (~\$2,400)

Metric	Mac Studio M4 Max 128GB	Dual RTX 3090 PC
Total memory	128 GB unified	48 GB VRAM (24+24)
Memory bandwidth	546 GB/s	~1,870 GB/s combined
Llama 70B Q4 speed	~15 tok/s (MLX)	~20-25 tok/s (vLLM)
Llama 32B Q4 speed	~28 tok/s	~60 tok/s
Power draw (load)	~120W	~700W
Power draw (idle)	~20W	~80W
Noise	Nearly silent	Loud under load
Physical size	7.7 x 7.7 x 3.7 inches	Mid-tower case
Total cost	~\$3,700	~\$2,400 (GPUs \$1,600 + system \$800)

The dual 3090 build is faster and cheaper. But it draws 6x the power, sounds like a jet engine under load, and can only fit models up to 48GB combined. If a 70B model at Q4 (~40GB) is your ceiling, the PC build wins on performance.

The Mac Studio wins when you need more than 48GB, when you care about noise, when it's sitting on a desk in a living room or office, or when you're running it 24/7 as an always-on server. Twenty watts at idle translates to about \$20/year in electricity. A dual 3090 rig idles at \$70-80/year and you'll hear the fans.

M3 Ultra 256GB (\$6,000) vs triple RTX 3090 PC (\$3,600)

Metric	Mac Studio M3 Ultra 256GB	3x RTX 3090 PC
Total memory	256 GB unified	72 GB VRAM
Memory bandwidth	819 GB/s	~2,800 GB/s combined
Llama 70B Q4 speed	~20 tok/s	~30+ tok/s
Can run 100B+ models	Yes	Needs 4-bit, barely fits
Power draw (load)	~215W	~1,100W
Total cost	~\$6,000	~\$3,600

Same pattern. The PC build is faster for models that fit. The Mac Studio loads models that don't fit on any reasonable PC configuration. The 256GB Ultra is for people who need to run very large models or multiple models simultaneously, and a triple-GPU rig can't match that memory pool.

The honest math

- **Tok/s per dollar:** NVIDIA wins. For the same budget, PC hardware generates tokens faster.
- **Memory per dollar:** Mac wins. 128GB of fast unified memory for \$3,700 has no PC equivalent.
- **Total cost of ownership:** Mac wins at 24/7 operation. Power savings add up over years.
- **Noise per tok/s:** Mac wins by a mile. The Mac Studio is inaudible under normal AI workloads. A multi-GPU PC is not.

Who should buy a Mac Studio for AI

Buy the M4 Max 128GB (\$3,699) if:

- You run 70B models regularly and need them loaded fast
- You want silent, always-on local AI serving on your network
- You build AI apps and need fast iteration on 14B-70B models
- You work in a shared space where GPU fan noise isn't acceptable
- You already use macOS and don't want to maintain a separate PC

Buy the M3 Ultra 256GB (\$6,000+) if:

- You're a researcher testing 70B+ models at high quantization

- You need multi-model setups (router + specialist architectures)
- You want to run the largest available open models (100B+ MoE)
- You fine-tune 14B+ models locally and keep hitting VRAM limits on GPUs

Who should NOT buy one

Don't buy a Mac Studio if:

- **14B models cover your needs.** A [MacBook Pro with 24GB](#) or even an [8GB Mac with the right models](#) handles that. The Mac Studio's value is memory headroom for larger models.
- **You need maximum tok/s.** A used RTX 3090 (\$700-900) plus a cheap PC generates tokens faster on models up to 24GB. Two of them beat the Mac Studio on everything up to 48GB.
- **You need CUDA.** PyTorch works on Metal, but some libraries, training frameworks, and inference tools only support NVIDIA. Check your toolchain before committing.
- **You're on a budget.** A [\\$500 budget AI PC](#) with an RTX 3060 12GB runs 7B-14B models fine. The Mac Studio is for people who've outgrown that tier.
- **Image generation is the priority.** NVIDIA's CUDA and Tensor cores still dominate Stable Diffusion, Flux, and ComfyUI. Mac Studio works for image gen, but it's slower per dollar than even mid-range NVIDIA cards.

Practical setup tips

Once you've decided, here's how to get the most out of it.

MLX vs Ollama

MLX is Apple's framework built for Apple Silicon. It's 20-30% faster than llama.cpp/Ollama for token generation. Use [MLX-LM](#) directly or LM Studio (which uses MLX as its backend on Mac).

If you want an API server, multi-model management, or integration with tools like Open WebUI, [Ollama is easier to set up](#). The speed difference vs MLX is noticeable but not dealbreaking for most workflows.

Keep models loaded

Set `OLLAMA_KEEP_ALIVE=-1`. The Mac Studio is meant to be always-on. Keep your primary model loaded in memory permanently. On a 128GB machine running a 40GB model, you've still got plenty of headroom for everything else.

Also set `OLLAMA_FLASH_ATTENTION=1` – it reduces memory usage with no quality loss. On a machine where you're pushing memory limits with large models, this can be the difference between fitting and not.

Use it as a server

The Mac Studio has 10Gb Ethernet, four Thunderbolt 5 ports, and draws 20W at idle. Set `OLLAMA_HOST=0.0.0.0`, enable Low Power Mode (new in 2025 models), and let it serve models to every device on your network. It's quieter and cheaper to run than any rack-mount server.

Watch memory pressure

Open Activity Monitor before loading your largest model. Green memory pressure means you're fine. Yellow is okay for occasional use. Red means the model is too large – drop to a lower quantization or smaller model. See our [Mac M-Series guide](#) for memory sizing details.

Bottom line

The Mac Studio is not the fastest AI hardware per dollar. It's the most memory per dollar in a silent, compact form factor. If your work requires models that exceed 24GB of VRAM, the Mac Studio solves that problem without a multi-GPU rig, without fan noise, and without a dedicated server room.

The M4 Max 128GB at \$3,699 is the configuration I'd recommend for most people doing serious local AI work. It runs 70B models, fits on a desk, and works as a general-purpose Mac when you're not running inference.

The M3 Ultra at 256GB+ is for researchers and developers who've already hit the ceiling of what 128GB can do and need to go bigger. It costs twice as much, but nothing else at any price gives you 256GB of fast unified memory in a box you can carry with one hand.

Get notified when we publish new guides.

[Subscribe – free, no spam](#)

Source: <https://insiderllm.com/guides/mac-studio-local-ai-workstation/>

Free guides for running AI locally