# Laptop vs Desktop for Local AI: Which Should You Buy?

January 31, 2026 · by Mark Bartlett

[Download this guide as PDF](#)

> **Quick Answer:** For raw performance per dollar, desktop wins decisively. A used RTX 3090 ($750) gives you 24GB VRAM — matching what costs $2,899+ in a laptop (RTX 5090 Laptop). The exception is MacBook Pro: a 48GB M4 Max ($3,999) can run 70B models that won't fit on any single desktop GPU under $2,000, thanks to unified memory. Gaming laptops with Windows are the worst option — the RTX 4070 Laptop only has 8GB VRAM despite costing $1,500+. Buy a desktop if you want the best value. Buy a MacBook if you need portability AND want to run large models. Avoid gaming laptops for AI unless you already own one.

📚 **More on this topic:** [Mac vs PC for Local AI](#) · [GPU Buying Guide](#) · [Budget AI PC Under $500](#) · [VRAM Requirements](#) · [Planning Tool](#)

The number one mistake people make when buying hardware for local AI: assuming a $2,000 gaming laptop will perform like a $2,000 desktop. It won't. Not even close.

A laptop RTX 4070 has 8GB VRAM. The desktop RTX 4070 has 12GB. A laptop RTX 4090 has 16GB. The desktop has 24GB. Same name, different chip, less memory. And for local AI, [VRAM is everything](#).

But there's a plot twist: MacBooks with Apple Silicon break the rules entirely. Their unified memory lets a $4,000 laptop run models that would need $3,000+ worth of NVIDIA desktop GPUs. So the answer isn't simply "buy a desktop." It depends on what you're actually doing.

## The Core Problem: VRAM per Dollar

For LLMs and image generation, the amount of memory your GPU can access determines what models you can run. Here's how desktop and laptop GPUs compare:

### RTX 40-Series: Desktop vs Laptop VRAM

| GPU Name | Desktop VRAM | Laptop VRAM | Desktop Price | Laptop System Price |
|----------|--------------|-------------|---------------|---------------------|
| RTX 4060 | 8 GB | 8 GB | ~$299 | ~$999-1,400 |

| GPU Name | Desktop VRAM | Laptop VRAM | Desktop Price | Laptop System Price |
|---|---|---|---|---|
| RTX 4070 | 12 GB | **8 GB** | ~$550 | ~$1,500-1,800 |
| RTX 4080 | 16 GB | **12 GB** | ~$800 (used) | ~$1,800-2,500 |
| RTX 4090 | 24 GB | **16 GB** | ~$1,600 | ~$2,500-4,000 |

## RTX 50-Series: Same Problem, New Generation

| GPU Name | Desktop VRAM | Laptop VRAM | Desktop Price | Laptop System Price |
|---|---|---|---|---|
| RTX 5060 | 8 GB | 8 GB | ~$299 | ~$1,099 |
| RTX 5070 | 12 GB | **8 GB** | ~$549 | ~$1,299 |
| RTX 5070 Ti | 16 GB | **12 GB** | ~$749 | ~$1,599 |
| RTX 5080 | 16 GB | 16 GB | ~$999 | ~$2,199 |
| RTX 5090 | 32 GB | **24 GB** | ~$1,999 | ~$2,899 |

The naming is deliberately misleading. An RTX 5090 Laptop uses a completely different GPU chip (GB203) than the desktop RTX 5090 (GB202). The desktop version has 107% more CUDA cores, 33% more VRAM, and 2x the memory bandwidth. They share a name and nothing else.

For local AI, this means a $1,500 gaming laptop with an RTX 4070 (8GB) can barely run a quantized 7B model with adequate context. A $750 used RTX 3090 in a desktop gives you 24GB — enough for 32B quantized models comfortably.

# Desktop Advantages

## More VRAM for Less Money

This is the biggest factor. Here's the cost per GB of usable AI memory:

| Hardware | AI Memory | Cost | Cost per GB |
|---|---|---|---|
| Used RTX 3090 (GPU only) | 24 GB | ~$750 | $31/GB |
| RTX 4060 desktop (GPU only) | 8 GB | ~$299 | $37/GB |
| RTX 5080 desktop (GPU only) | 16 GB | ~$999 | $62/GB |

| Hardware | AI Memory | Cost | Cost per GB |
|---|---|---|---|
| RTX 5090 Laptop (full system) | 24 GB | ~$2,899 | $121/GB |
| RTX 4090 Laptop (full system) | 16 GB | ~$2,500-4,000 | $156-250/GB |

Even when you add the cost of a full desktop system ($500-800 for CPU, motherboard, RAM, PSU, case, and SSD), a desktop with a used RTX 3090 comes out to ~$1,250-1,550 total for 24GB VRAM. The cheapest laptop with equivalent VRAM (RTX 5090 Laptop) starts at $2,899.

### Upgradable

Buy a desktop with a mid-range GPU now, upgrade later. Your CPU, RAM, case, and PSU carry forward. You can't upgrade a laptop GPU — what you buy is what you're stuck with.

### Better Sustained Performance

Desktop GPUs run at full power (200-575W) with large heatsinks and case fans. Laptop GPUs run at 50-175W with thin cooling solutions. Under sustained AI workloads — which keep the GPU at 100% utilization with no breaks — laptops throttle.

Real-world impact: one documented test showed a laptop dropping from 12.4 tok/s to 4.1 tok/s during sustained LLM inference — a 67% performance loss from thermal throttling. Desktop GPUs with adequate cooling don't have this problem.

### Used Market Access

You can buy a used RTX 3090 for ~$750. You can't buy a used laptop GPU and install it. Used gaming laptops exist, but they depreciate fast, have worn batteries, and you can't verify thermal paste condition.

## Laptop Advantages

### Portability

The obvious one. If you need AI on the go — demos, travel, working from coffee shops — a desktop isn't an option. And some people genuinely need portability more than raw performance.

### All-in-One

A laptop includes screen, keyboard, trackpad, battery, speakers, webcam, and WiFi. A desktop equivalent costs more when you factor in a monitor, peripherals, and desk space.

### MacBooks Break the Rules

This is the big one. MacBook Pros with Apple Silicon aren't bound by the VRAM limitation that cripples Windows gaming laptops. More on this below.

## MacBook Pro: The Laptop Exception

Apple Silicon uses unified memory — the CPU, GPU, and Neural Engine all share the same memory pool at full bandwidth. There's no 8GB or 16GB VRAM ceiling. A 48GB MacBook Pro can load a 48GB model. A 128GB MacBook Pro can load a 128GB model.

This changes the math completely for local AI.

### What Each MacBook Config Can Run

| Config | Unified Memory | Bandwidth | What You Can Run (Q4 quantized) | Price |
|---|---|---|---|---|
| M4 Pro | 24 GB | 273 GB/s | 14B comfortably, 27B tight | $1,999 |
| M4 Pro | 48 GB | 273 GB/s | 32B comfortably, 70B tight | $2,899 |
| M4 Max (32c GPU) | 36 GB | 410 GB/s | 27B comfortably | $3,199 |
| M4 Max (40c GPU) | 48 GB | 546 GB/s | 32B comfortably, 70B tight | $3,999 |
| M4 Max (40c GPU) | 128 GB | 546 GB/s | 70B comfortably, 100B+ possible | ~$5,199 |

### Mac vs NVIDIA Speed Comparison

| Model | M4 Max 40c (546 GB/s) | RTX 4090 (1,008 GB/s) |
|---|---|---|
| 8B Q4 | ~83 tok/s | ~130 tok/s |
| 14B Q4 | ~40-50 tok/s | ~60-80 tok/s |
| 70B Q4 | ~12 tok/s | **2-5 tok/s** (doesn't fit in 24GB, offloads to CPU) |

That last row is the key. An RTX 4090 is faster for any model that fits in 24GB VRAM. But a 70B model at Q4 needs ~37GB — it doesn't fit, so the NVIDIA card has to offload to system RAM over PCIe, and performance collapses. The Mac runs it at full unified memory bandwidth.

## Mac Software for Local AI

- **Ollama**: Works out of the box, auto-detects Metal. Easy but not the fastest option (~20-40 tok/s for small models).
- **LM Studio**: Supports both GGUF (llama.cpp) and MLX models. MLX backend is 20-50% faster than Ollama for the same models. The best GUI option on Mac.
- **MLX (command line)**: Apple's own ML framework, optimized for Apple Silicon. Fastest option — up to 230 tok/s for 7B models on M2 Ultra, ~83 tok/s for 8B on M4 Max.
- **ComfyUI**: Works for Stable Diffusion and Flux image generation via PyTorch MPS, but 3-5x slower than NVIDIA CUDA.
- **mflux**: MLX-native Flux implementation. Generates 1024x1024 Flux Schnell images in ~10.5 seconds on M4 Max — competitive with RTX 4090.

For a deeper comparison, see the Mac vs PC for Local AI guide.

## When to Buy a MacBook for AI

- You need portability AND want to run models larger than 32B parameters
- You value silence and efficiency (40-80W total system power vs 300-575W for desktop GPU alone)
- You're willing to pay a premium for the unified memory advantage
- You'll use it for other work too (development, creative work, daily computing)

## When NOT to Buy a MacBook for AI

- You only run 7B-14B models (cheaper to build a desktop with an RTX 3060 12GB or 4060 for $800)
- You do heavy image generation (NVIDIA is 3-5x faster for Stable Diffusion/Flux via ComfyUI)
- You want to fine-tune models (CUDA ecosystem is far ahead for training)
- Budget is the primary concern

# Gaming Laptops: The Worst Value for AI

Windows gaming laptops are the worst value proposition for local AI. Here's why:

## The VRAM Problem

Most gaming laptops in the $1,000-2,000 range have 8GB VRAM. That's it. The RTX 4060 Laptop, RTX 4070 Laptop, RTX 5060 Laptop, RTX 5070 Laptop — all 8GB. For local AI, 8GB limits you to 7B quantized models with tight context windows.

To get more than 8GB VRAM in a Windows laptop, you need:

| Laptop VRAM | Minimum Laptop Price | Desktop Equivalent Cost |
|---|---|---|
| 12 GB | ~$1,599 (RTX 5070 Ti Laptop) | ~$250 (used RTX 3060 12GB) |
| 16 GB | ~$2,199 (RTX 5080 Laptop) | ~$720 (used RTX 4070 Ti Super) |
| 24 GB | ~$2,899 (RTX 5090 Laptop) | ~$750 (used RTX 3090) |

A used RTX 3090 desktop card costs $750 and gives you 24GB VRAM with 936 GB/s bandwidth. The cheapest laptop with 24GB VRAM costs $2,899 and gives you lower bandwidth through a different chip. The desktop card alone costs 74% less.

## Thermal Throttling Is Real

AI workloads keep the GPU at 100% utilization continuously — unlike gaming, which fluctuates with scene complexity. Laptop cooling isn't designed for this:

- Sustained temperatures of 78-88°C are normal
- Thermal throttling can cut throughput by 40-67%
- Long-term heat stress degrades components faster
- Fan noise is substantial during sustained loads

**Mitigation helps but doesn't solve the problem.** Undervolting, cooling pads, and smaller quantizations reduce throttling but you're still fighting the form factor. A desktop with a $30 tower cooler doesn't have these issues.

## If You Already Own a Gaming Laptop

Don't buy a new one for AI. Instead:

- Use it for small models (7B-8B) — this works fine on 8GB VRAM

- For larger models, either build a budget desktop (under $500 is possible) or use cloud APIs for the occasional heavy task
- Run CPU-only inference for models up to 3-4B parameters if you have 16GB+ system RAM

## CPU-Only on Laptops

If your laptop has no dedicated GPU (or only a weak one), you can still run models on CPU — just slowly.

| CPU Type | 7B Model (Q4) | 1-3B Model (Q4) |
|---|---|---|
| Modern x86 (Intel 13th/14th gen, AMD 7000/8000) | ~5-15 tok/s | ~40-50 tok/s |
| Apple M4 | ~15-25 tok/s | ~50-80 tok/s |
| Apple M4 Pro | ~25-40 tok/s | ~70-100 tok/s |
| Older x86 (Intel 10th/11th gen) | ~3-8 tok/s | ~20-30 tok/s |

On x86 laptops, memory bandwidth is the bottleneck. DDR5 helps over DDR4, and dual-channel matters. For the best small models to run on CPU, see our dedicated guide.

Apple Silicon is significantly faster at CPU-only inference because its unified memory bandwidth (120-546 GB/s) is much higher than typical laptop DDR5 (~50-70 GB/s). Even the base M4 MacBook Air outperforms most x86 laptops with discrete GPUs for models that fit in memory.

## The Decision Matrix

### Budget: Under $1,000

**Buy: Desktop.** A budget AI PC under $500 with a used RTX 3060 12GB ($200) runs 7B-14B models at good speeds. No laptop at this price comes close.

### Mid-Range: $1,000-2,000

**Buy: Desktop.** A used RTX 3090 ($750) in a $1,500 total desktop build gives you 24GB VRAM — enough for 32B models. The best laptop you can get at this price has 8GB VRAM.

## Premium: $2,000-3,000

**Buy: Desktop or MacBook Pro.** A desktop with an RTX 5080 ($999, 16GB) outperforms any Windows laptop. A MacBook Pro M4 Pro with 48GB ($2,899) runs larger models thanks to unified memory but is slower per-token for models that fit in NVIDIA VRAM.

## High-End: $3,000+

**Buy: MacBook Pro M4 Max or Desktop RTX 5090.**

- If you need portability: MacBook Pro M4 Max with 64-128GB unified memory. Can run 70B+ models that no single consumer NVIDIA GPU handles. Slower per-token but runs things nothing else can.
- If you don't need portability: Desktop RTX 5090 (32GB, $1,999) in a $3,500 total build. Fastest single-GPU option for models up to 32B. For 70B+, you'd need dual GPUs.

## Already Own a Laptop

**Don't replace it.** Use it for what it can handle (small models, CPU inference, cloud APIs). If you want more AI capability, add a desktop instead of upgrading the laptop. You keep the laptop for portability and the desktop for heavy lifting.

---

# Can You Use Both? eGPU and Hybrid Setups

### External GPU (eGPU) via Thunderbolt

You can connect a desktop GPU to a laptop through a Thunderbolt enclosure. The reality:

| Connection | Bandwidth | Performance Hit vs Desktop |
|---|---|---|
| PCIe 4.0 x16 (desktop) | ~32 GB/s | 0% (reference) |
| Thunderbolt 3/4 / USB4 | ~5 GB/s | ~38% slower |
| Thunderbolt 5 | ~10-15 GB/s | ~20% slower (estimated) |
| OCuLink (PCIe 4.0 x4) | ~8 GB/s | ~25% slower |

A Thunderbolt eGPU with an RTX 3090 tested at 38.5% lower tok/s compared to the same card on PCIe. That's a significant penalty, but you still get access to 24GB VRAM — which matters more than raw speed for many use cases.

**Worth it if:** You already own a laptop and want desktop-class VRAM without building a full desktop. The eGPU enclosure costs $200-400, plus the GPU.

**Not worth it if:** You're buying from scratch. Building a full desktop costs less and performs better.

### Hybrid Approach

The most practical setup for many people: use a laptop (especially MacBook) for small models, daily chat, and portable use. Keep a desktop with a high-VRAM GPU for larger models, image generation, and heavy workloads. Remote access via SSH or a web UI like Open WebUI lets you use the desktop's GPU from anywhere on your home network.

## Bottom Line

**Desktop is the best value for local AI.** Nothing else comes close on VRAM per dollar. A $1,500 desktop build with a used RTX 3090 (24GB) outperforms a $3,000 gaming laptop with an RTX 5090 Laptop (24GB VRAM but slower chip and thermal throttling).

**MacBook Pro is the best laptop for local AI.** Unified memory changes the equation — a 48-128GB MacBook can run models that exceed any single desktop GPU's VRAM. It's slower per-token for smaller models, but it can run things nothing else can in a portable form factor.

**Gaming laptops are the worst option for local AI.** 8GB VRAM at $1,000-2,000 is a terrible deal when a desktop RTX 3060 12GB costs $200 used. Only consider a gaming laptop if you already own one and don't want to buy new hardware.

Build a desktop if you can. Buy a MacBook if you need portability. Use what you already have if it's "good enough." And whatever you do, don't buy a $2,000 gaming laptop expecting it to be a local AI workstation.

Get notified when we publish new guides.

Subscribe — free, no spam

Source: https://insiderllm.com/guides/laptop-vs-desktop-local-ai/

Free guides for running AI locally