# Apple M5 Pro and M5 Max: What 4x Faster LLM Processing Actually Means for Local AI

March 3, 2026 · by Mark Bartlett

[Download this guide as PDF](#)

> **Quick Answer:** The M5 Pro (307GB/s, up to 64GB) and M5 Max (614GB/s, up to 128GB) are the largest Apple Silicon upgrades for local AI since M1. The 4x faster LLM prompt processing comes from Neural Accelerators embedded in every GPU core, replacing the old fixed-size Neural Engine as the main AI workhorse. Token generation speed scales with bandwidth, and the M5 Max is 12% faster than the M4 Max (614 vs 546GB/s). The M5 Max 128GB is now the best single device for running 70B+ models portably. Pre-orders March 4, ships March 11. Worth upgrading from M1/M2 Pro/Max, harder to justify from M4 Pro/Max unless you need the prompt processing speed.

📚 **More on this topic:** [Best Local LLMs for Mac](#) · [Running LLMs on Mac M-Series](#) · [Mac vs PC for Local AI](#) · [VRAM Requirements](#)

Apple just announced the M5 Pro and M5 Max MacBook Pros, and I think they buried the lead in the spec sheet. Forget the CPU speed bump. The M5 Max hits 614GB/s memory bandwidth. Neural Accelerators are built into every GPU core. And Apple explicitly name-dropped [LM Studio](#) in the official press release.

They're marketing MacBook Pro as an LLM machine now. Here's what actually matters for running models locally.

---

## What Apple Announced

Pre-orders open March 4. Ships March 11, 2026.

| Spec | M5 Pro | M5 Max |
|---|---|---|
| **CPU** | 18-core (6 super + 12 performance) | 18-core (6 super + 12 performance) |
| **GPU** | Up to 20-core | Up to 40-core |
| **Neural Accelerators** | 1 per GPU core (up to 20) | 1 per GPU core (up to 40) |
| **Neural Engine** | 16-core | 16-core |

| Spec | M5 Pro | M5 Max |
|---|---|---|
| Unified Memory | Up to 64GB | Up to 128GB |
| Memory Bandwidth | 307GB/s | 614GB/s |
| Process | Dual 3nm dies (Fusion Architecture) | Dual 3nm dies (Fusion Architecture) |
| SSD | Up to 14.5GB/s read/write | Up to 14.5GB/s read/write |
| Starting Storage | 1TB | 2TB |

The headline claims: 4x faster LLM prompt processing vs M4 Pro/Max, up to 8x faster AI image generation vs M1 Pro/Max, and 30% faster CPU performance.

## Pricing

| Configuration | Price |
|---|---|
| 14" M5 Pro | $2,199 |
| 16" M5 Pro | $2,699 |
| 14" M5 Max | $3,599 |
| 16" M5 Max | $3,899 |

The 128GB M5 Max configs will run $5,000+ once you add storage upgrades. That's relevant when we get to the "is it worth it" section.

# The Two Numbers That Matter for Local AI

If you run LLMs locally, two specs determine your experience: memory capacity (how big a model you can load) and memory bandwidth (how fast it generates tokens).

## Memory Bandwidth: M5 Max Jumps to 614GB/s

Here's how M5 stacks up against every previous generation:

| Chip | Memory Bandwidth | Relative Speed |
|---|---|---|
| M1 Pro | 200 GB/s | 1x |
| M1 Max | 400 GB/s | 2x |

| Chip | Memory Bandwidth | Relative Speed |
|---|---|---|
| M2 Pro | 200 GB/s | 1x |
| M2 Max | 400 GB/s | 2x |
| M3 Pro | 150 GB/s | 0.75x |
| M3 Max | 300-400 GB/s | 1.5-2x |
| M4 Pro | 273 GB/s | 1.4x |
| M4 Max | 546 GB/s | 2.7x |
| **M5 Pro** | **307 GB/s** | **1.5x** |
| **M5 Max** | **614 GB/s** | **3.1x** |

Token generation for LLMs is memory-bandwidth-bound. Apple's own MLX research confirms this: the base M5 generates tokens 19-27% faster than the M4, matching its 28% bandwidth increase (153 vs 120GB/s). The relationship is nearly linear.

Apply that math to the Pro and Max tiers. The M5 Pro at 307GB/s should generate tokens about 12% faster than the M4 Pro at 273GB/s. The M5 Max at 614GB/s should be about 12% faster than the M4 Max at 546GB/s.

That's not 4x. So where does the 4x number come from?

## The 4x Claim: Prompt Processing, Not Token Generation

Apple says "4x faster LLM prompt processing." That's prompt processing, specifically the time-to-first-token (TTFT). This is the compute-bound phase where the model digests your entire input before generating the first response token.

The MLX benchmarks on the base M5 back this up. Time-to-first-token improved 3.3-4x across models:

| Model | TTFT Speedup (M5 vs M4) |
|---|---|
| Qwen 1.7B (BF16) | 3.57x |
| Qwen 8B (BF16) | 3.62x |
| Qwen 8B (4-bit) | 3.97x |
| Qwen 14B (4-bit) | 4.06x |
| GPT OSS 20B (MXFP4) | 3.33x |

| Model | TTFT Speedup (M5 vs M4) |
|---|---|
| Qwen 30B MoE (4-bit) | 3.52x |

These are base M5 numbers. The Pro and Max should hit similar or higher speedups for prompt processing because they have more Neural Accelerators (20 and 40 respectively vs 10 on the base M5).

What does faster prompt processing mean in practice? If you're feeding a large document into a model for summarization, or doing RAG with big context windows, the wait before the model starts responding drops dramatically. On an M4 Pro, a 32K-token prompt on a 14B model might take 8-10 seconds to process. On an M5 Pro, that drops to 2-3 seconds.

But once generation starts, you're bandwidth-limited. The M5 Pro will generate tokens roughly 12% faster than the M4 Pro. Still an improvement, but not 4x.

## Neural Accelerators: The Actual Architecture Change

Previous Apple Silicon had a fixed 16-core Neural Engine for AI workloads. Same 16 cores whether you bought the base chip or the Max. The M5 changes this: every GPU core now has its own Neural Accelerator built in. The M5 Max with 40 GPU cores gets 40 Neural Accelerators. The 16-core Neural Engine is still there too.

This matters more than the spec sheet suggests.

AI compute now scales with GPU core count. The M5 Max gets 40 Neural Accelerators vs 20 on the M5 Pro. When Apple releases the M5 Ultra (which will fuse two M5 Max dies), that's 80 Neural Accelerators. Previously, every tier was stuck with the same 16-core Neural Engine.

It also means MLX workloads benefit directly. MLX runs on Metal, which uses GPU cores. Neural Accelerators embedded in those same cores means MLX can use them without shuffling data to a separate chip. The TTFT benchmarks back this up: the 3.5-4x speedup on compute-bound prompt processing comes from these Neural Accelerators doing matrix multiplications alongside the GPU shader cores.

## What Can You Actually Run on M5 Pro and M5 Max

Here's the practical breakdown, based on model sizes and the memory each config provides:

| Configuration | Unified Memory | Bandwidth | Models You Can Run | Expected Token Speed |
|---|---|---|---|---|
| **M5 Pro 36GB** | 36GB | 307 GB/s | 30B Q4, 14B Q8, all 7-8B models | 14B: ~35-45 tok/s, 30B Q4: ~18-25 tok/s |
| **M5 Pro 64GB** | 64GB | 307 GB/s | 70B Q4, 30B Q8, multiple models loaded | 70B Q4: ~12-18 tok/s, 30B Q8: ~18-25 tok/s |
| **M5 Max 64GB** | 64GB | 614 GB/s | Same models, ~2x faster generation | 70B Q4: ~22-32 tok/s, 30B Q8: ~35-45 tok/s |
| **M5 Max 128GB** | 128GB | 614 GB/s | 120B+ Q4, 70B Q8, frontier models | 70B Q8: ~20-28 tok/s, 120B Q4: ~12-16 tok/s |

The M5 Max 128GB is the config I'd be most excited about. You could load Llama 3.3 70B at Q8 (full quality, ~75GB) with plenty of room for context and system memory, at speeds that are actually usable. On PC, running 70B at Q8 requires either a 48GB GPU (A6000, $4,000+) or multi-GPU setups.

For reference, the M5 Max 128GB at 614GB/s should generate roughly 22-32 tok/s on 70B Q4 models via MLX. That's competitive with an RTX 4090 running the same model in GGUF (which requires CPU offloading at 70B, killing performance).

## MLX vs Ollama on M5: Which to Use

Both work on M5, but there's a clear performance gap.

MLX is Apple's own ML framework, optimized for Metal (and now Metal 4 on M5). It benefits most from Neural Accelerators because it runs directly on GPU cores where those accelerators live. On the base M5, MLX already shows the 3.5-4x TTFT improvement. Expect similar gains on Pro and Max.

Ollama uses llama.cpp under the hood, which has Metal support but isn't specifically optimized for Neural Accelerators. It'll still benefit from the increased bandwidth (faster token generation) but probably won't fully exploit the TTFT improvements until llama.cpp adds Neural Accelerator support.

LM Studio now offers an MLX backend on Mac, which gives you MLX performance with a GUI. Apple mentioning LM Studio by name in the press release wasn't an accident.

My recommendation: use MLX (directly or via LM Studio's MLX backend) to get the full benefit of M5 hardware. Ollama is fine for convenience, but you're leaving prompt processing performance on the table. See our Ollama vs LM Studio comparison for a deeper breakdown.

## Should You Upgrade?

This depends entirely on what you're coming from and whether you run local models regularly.

### From M1 Pro/Max: Yes

A 6-8x AI performance jump. The M1 Max at 400GB/s was solid for its time, but the M5 Max at 614GB/s with Neural Accelerators in every GPU core is a different tier. If you're running local models regularly, this is worth the upgrade. Your five-year-old machine has had a good run.

### From M2 Pro/Max: Meaningful

The bandwidth jump is real (200GB/s to 307GB/s on Pro, 400GB/s to 614GB/s on Max), and the Neural Accelerator architecture is genuinely new. If you're hitting memory limits on a 32GB M2 Max, the 128GB M5 Max opens up model sizes you couldn't touch before. Not urgent, but worth considering if you're running models daily.

### From M3 Pro/Max: Moderate

The M3 Max at 300-400GB/s to M5 Max at 614GB/s is a real bandwidth jump, and the Neural Accelerator architecture is genuinely new. Prompt processing will be 3-4x faster. But if your M3 Max handles your current workloads fine, you could wait for the M5 Ultra if you want the maximum possible.

### From M4 Pro/Max: Hard Sell

The M4 Max at 546GB/s to M5 Max at 614GB/s is only a 12% bandwidth increase. Token generation speed won't feel dramatically different. The 4x prompt processing improvement is real and you'll notice it with large contexts, but the M4 Max 128GB is still an excellent machine for local AI. Wait for the M5 Ultra unless fast prompt processing is a specific pain point.

### From Any Intel Mac

Stop reading this section and go order one.

## Price vs Alternatives: Honest Comparison

A maxed-out M5 Max MacBook Pro with 128GB will cost somewhere in the $5,000-7,000 range once you add storage. That's a lot of money. Let's compare alternatives honestly.

| Platform | Config | Cost | Bandwidth | 70B Q4 Speed | Portable? |
|---|---|---|---|---|---|
| **M5 Max MBP** | 128GB unified | ~$5,000-7,000 | 614 GB/s | ~22-32 tok/s | Yes |
| **Desktop + RTX 4090** | 24GB VRAM + 64GB RAM | ~$2,500 | 1,008 GB/s (VRAM) | ~25-35 tok/s (needs offload) | No |
| **Desktop + RTX 5090** | 32GB VRAM + 64GB RAM | ~$3,500 | 1,792 GB/s (VRAM) | ~40-55 tok/s (needs offload) | No |
| **AMD Strix Halo Mini PC** | 128GB unified | ~$1,700-2,200 | ~273 GB/s | ~10-15 tok/s | Small form factor |

The Mac wins on portability, silence, and battery life (up to 22 hours). And it can load a 70B model entirely in unified memory without offloading. The RTX 4090 is faster per dollar but can't fit 70B in 24GB VRAM, so it has to offload to system RAM, which tanks performance.

The Mac loses on raw price-to-performance. A $2,500 desktop with an RTX 4090 is faster for models that fit in 24GB. An AMD Strix Halo mini PC with 128GB costs a third of the price with the same memory capacity, though at less than half the bandwidth.

Where the Mac has no competition: 614GB/s bandwidth and 128GB unified memory in a laptop. If you need to run 70B+ models while traveling, there's no alternative right now.

## Fusion Architecture and the M5 Ultra Question

The M5 Pro and M5 Max are each built from two third-generation 3nm dies fused into a single SoC via a high-bandwidth interconnect. Apple calls this "Fusion Architecture." It's similar to what AMD does with chiplets, but Apple's interconnect is fast enough that the two dies behave as one chip from a software perspective.

Here's where it gets interesting for the wait-and-see crowd. Apple has fused two Max dies to make every Ultra chip since M1. If they do it again, the M5 Ultra would have up to 80 GPU cores, 80 Neural Accelerators, 256GB unified memory, and roughly 1,228GB/s bandwidth. That would make the Mac Studio the most capable single device for running frontier-class models (200B+) locally.

No official M5 Ultra details yet. But if you can wait and you need maximum capacity, it's almost certainly coming later this year.

## What This Means for Local AI

Apple putting LM Studio in an official press release is a statement. "Run advanced LLMs on device" and "train custom models locally" in the marketing copy. A year ago, Apple's AI messaging was all about Siri and on-device photo features. Now they're talking about token generation speeds and model training. That's a shift worth paying attention to.

The Neural Accelerator architecture is what I'd focus on more than the 4x headline. AI compute now scales with GPU core count, which means every future Apple Silicon chip gets proportionally better at LLM inference as Apple adds GPU cores. The old fixed Neural Engine was a ceiling. This removes it.

And by embedding those accelerators in GPU cores that MLX targets directly, Apple is making their own framework the fastest path for local inference on Mac. Third-party tools will catch up, but MLX has a structural advantage on M5 hardware.

For those of us who've been running models on Mac for a while: this is the biggest upgrade cycle since the original M1 made Apple Silicon viable for AI work. The bandwidth jump alone changes which models you can comfortably run, and the Neural Accelerator architecture means the M5 will age better than previous generations as MLX and other tools optimize for it.

## Bottom Line

The M5 Pro and M5 Max are a genuine architecture change for local AI. The 4x prompt processing claim is real for compute-bound workloads, powered by Neural Accelerators in every GPU core. Token generation is 12% faster, tracking the bandwidth increase. The M5 Max 128GB at 614GB/s is the best portable device for running large language models right now.

If you're on M1 or M2, upgrade. If you're on M4 Max, you're probably fine unless fast prompt processing on large contexts is worth $4,000+ to you. If you're on Intel, you know what to do.

Pre-orders open March 4 at 6:15 a.m. PT. Ships March 11.

📚 **Related guides:** Best Local LLMs for Mac 2026 · Running LLMs on Mac M-Series · Mac vs PC for Local AI · Stable Diffusion on Mac with MLX · GPU Buying Guide · VRAM Requirements

Get notified when we publish new guides.

Subscribe — free, no spam

---

Source: https://insiderllm.com/guides/apple-m5-pro-max-local-ai/

Free guides for running AI locally