# Engineering High-Concurrency Inference Systems with Apple MLX: A Comprehensive Technical Analysis (2026)

## 1. Executive Summary

As of January 2026, the ecosystem surrounding Apple's MLX framework has matured from a promising research experiment into a robust substrate for high-performance machine learning on Apple Silicon. The introduction of the M5 chip series, with its enhanced Neural Accelerators and massive memory bandwidth, has shifted the bottleneck of local inference from raw compute capability to system architecture and memory orchestration. For machine learning engineers and systems architects, the primary challenge is no longer merely executing a Large Language Model (LLM) on a Mac, but serving it efficiently under concurrent load.

This report provides an exhaustive technical analysis of achieving concurrent query processing using MLX. It posits that while the platform's Unified Memory Architecture (UMA) offers a distinct theoretical advantage over traditional discrete GPU setups—specifically by eliminating the PCIe transfer tax—realizing this potential requires a radical departure from standard CUDA-centric design patterns. True concurrency in MLX is not achieved through simple parallelism but through **Continuous Batching**, rigorous **Key-Value (KV) cache management**, and a sophisticated interplay between Python's asynchronous I/O and MLX's lazy evaluation graph.

The analysis delineates the critical boundary between the conveniences provided by the mlx-lm package and the custom engineering required for production-grade serving. While mlx-lm offers optimized kernels and basic batching primitives, high-throughput systems require custom schedulers, memory managers, and compilation strategies. This document serves as a blueprint for bridging that gap, offering a detailed exploration of optimization vectors ranging from quantization-aware batch sizing to compiler-friendly tensor padding.

## 2. The Architectural Foundation: MLX on Apple Silicon

To engineer effective concurrent systems, one must first deconstruct the underlying interaction between the MLX software stack and the Apple Silicon hardware. The concurrency model in MLX is fundamentally different from frameworks like PyTorch on NVIDIA GPUs, driven by two unique characteristics: Unified Memory and Lazy Evaluation.

### 2.1 The Unified Memory Advantage in Concurrency

In traditional server-grade architectures involving discrete GPUs (dGPUs), memory is bifurcated. The CPU has its host memory (DRAM), and the GPU has its device memory (VRAM). Concurrent inference on such systems involves a complex dance of moving data across the PCIe bus. When a new batch of requests arrives, the input tokens must be copied to VRAM, and generated tokens must be copied back. This transfer latency puts a hard floor on response times, particularly for high-frequency, low-latency requests.

Apple Silicon's Unified Memory Architecture (UMA) fundamentally alters this equation. In an M4 or M5 powered Mac Studio, the CPU and GPU share a single pool of high-bandwidth memory. An mx.array created by the CPU is effectively a pointer to a memory address that the GPU can read immediately without a copy operation.

**Implications for Concurrency:**

- **Zero-Copy Batch Construction:** The CPU can assemble a batch of concurrent user prompts (performing tokenization, padding, and layout organization) directly in the memory space that the GPU will consume. This eliminates the "prefill latency" associated with data transfers.[1]
- **Oversubscription Viability:** Because there is no hard VRAM wall, the system can more gracefully handle bursts of concurrent requests that might momentarily exceed the "ideal" working set, leveraging the OS's compressed memory and paging mechanisms, although performance degrades, it does not crash as readily as OOM (Out of Memory) errors on dGPUs.

## 2.2 The GIL, Threads, and Lazy Evaluation

A persistent myth in the Python ecosystem is that the Global Interpreter Lock (GIL) renders Python unsuitable for concurrent threading. In the context of MLX, this limitation is largely neutralized.

MLX operations are **lazy**. When a Python thread executes a command like c = mx.matmul(a, b), the framework does not immediately perform the matrix multiplication. Instead, it records an operation node in a computation graph and returns a handle to the future result. This graph construction is extremely fast and CPU-light. The actual heavy computation is triggered only when mx.eval(c) is called or the data is implicitly required (e.g., for printing or conditional logic).[2]

Crucially, when the computation is triggered, MLX releases the GIL. The heavy lifting is offloaded to the Metal API, running on the GPU or Neural Engine. During this time, the Python interpreter is free to switch to another thread.

**The Concurrency Cycle:**

1. **Thread A (Request 1):** Prepares input, calls model(input), and hits mx.eval(). MLX releases the GIL and dispatches work to the GPU.
2. **Thread B (Request 2):** Wakes up immediately. It prepares its input, calls model(input), and hits mx.eval(). MLX dispatches this work to the command queue.
3. **Thread C (Network):** Handles an incoming HTTP request, parses JSON, and places it in a queue.

This behavior transforms LLM inference from a CPU-bound task (which Python handles

poorly) to an I/O-bound-like task (which Python handles well). Consequently, standard Python threading or asyncio patterns become highly effective for orchestrating concurrent MLX workloads.[4]

## 2.3 Streams and Hardware Scheduling

While Python threads allow for concurrent *dispatch*, the hardware execution is serialized unless explicitly managed. MLX uses the concept of **streams** (similar to CUDA streams). By default, all operations run on the default stream of the default device.

If multiple Python threads dispatch operations to the same default stream, the GPU will execute them sequentially. For LLM inference, this is often desirable because the model weights are huge; running two inference passes purely in parallel would trash the memory cache. The goal of concurrency in LLMs is not necessarily to run two matrix multiplications at the exact same nanosecond, but to maximize the **throughput** of the pipeline so that the GPU is never idle.

However, for auxiliary tasks—such as decoding the next token for User A while simultaneously pre-processing the prompt for User B—using separate streams can unlock genuine hardware parallelism, allowing the CPU/Neural Engine to work alongside the GPU.[6]

# 3. Achieving Concurrent Queries: Implementation Patterns

The user's query explicitly asks *how* to achieve concurrent queries. The answer lies in moving up the stack from raw hardware capabilities to software architecture. As of 2026, the industry standard has converged on the **Asyncio Event Loop** pattern driving a **Continuous Batching** engine.

## 3.1 The "Naive" Approach vs. Production Reality

A naive implementation might involve spinning up a FastAPI server where every endpoint simply calls model.generate().
- **The Problem:** model.generate() is typically a blocking operation that runs until the entire sequence is finished. If User A asks for a 500-word essay, User B's request will hang until User A is done. This is **Head-of-Line (HOL) blocking**.
- **The Threading Trap:** Spawning a new thread for each request does not solve this if they all try to access the model simultaneously without batching. They will simply fight for the GPU, leading to thrashing and high latency for everyone.

## 3.2 The Asyncio Queue Architecture

To solve HOL blocking, we decouple the *reception* of a request from its *execution*.
**The Architecture:**
1. **Ingestion Layer (asyncio):** A lightweight asynchronous web server (FastAPI/Uvicorn) receives requests. It validates inputs and pushes a RequestObject (containing the prompt, parameters, and a Future for the result) into a global asyncio.Queue.

2.  **Scheduler Layer (Background Task):** A dedicated "Runner" coroutine sits in an infinite loop. It pulls requests from the queue. Crucially, it does not process them one by one. It implements a **scheduling policy** (see Section 3.3) to group multiple requests into a single batch.
3.  **Execution Layer (MLX):** The scheduler passes the batch to the model. The model performs *one single step* of inference (generating one token for each request in the batch).
4.  **Distribution Layer:** The generated tokens are sent back to the respective RequestObjects. Completed requests are removed from the batch; active ones remain for the next step.

This architecture ensures that the server remains responsive. Even if the GPU is 100% utilized, the network layer can accept new connections and queue them.[8]

## 3.3 Continuous Batching (Iteration-Level Scheduling)

Static batching (waiting for $N$ requests to arrive before starting) is obsolete. The dominant technique in 2026 is **Continuous Batching** (also known as cellular batching or iteration-level scheduling).

**Mechanism:**
Imagine a bus (the GPU batch) that never stops.

-  **Step 1:** The bus has 3 passengers (Requests A, B, C). It drives one mile (generates one token).
-  **Step 2:** Passenger B reaches their destination (generates an EOS token). They get off (are removed from the batch).
-  **Step 3:** The bus now has an empty seat. Passenger D is waiting at the stop (Queue). Passenger D gets on immediately.
-  **Step 4:** The bus drives another mile (generates the next token for A, C, and the first token for D).

**Implementation in MLX:**
This logic is *not* entirely encapsulated in mlx_lm.generate. While mlx-lm supports batch input, it assumes a static list. To achieve continuous batching, one must write a custom inference loop using mlx_lm.models.base.generate_step.

**Code Concept:**

Python

```
# Conceptual logic for Continuous Batching in MLX
active_requests =
while True:
    # 1. New Arrivals: Check queue non-blocking
    while not queue.empty() and len(active_requests) < MAX_BATCH_SIZE:
        req = queue.get_nowait()
        active_requests.append(req)
```

```
# 2. Prune Completed: Remove finished requests
active_requests = [r for r in active_requests if not r.is_finished]

if not active_requests:
    await asyncio.sleep(0.001) # Avoid busy loop
    continue

# 3. Form Batch: Collect inputs from all active requests
# Note: Some are just starting (prompts), others are continuing (tokens)
batch_inputs = prepare_batch(active_requests)

# 4. Inference Step: Run the model
# The 'step' function handles the math
next_tokens, new_cache = model.generate_step(batch_inputs, cache)

# 5. Distribute Results
for i, req in enumerate(active_requests):
    req.stream.send(next_tokens[i])
```

This loop maximizes hardware utilization (High Throughput) while minimizing the wait time for new requests (Low Latency).[8]

# 4. Optimization Strategy: Maximizing Hardware Efficiency

Once the concurrent architecture is in place, the focus shifts to optimization. In MLX, optimization is a function of memory bandwidth, compilation, and precision.

## 4.1 Quantization: The Bandwidth Multiplier

On Apple Silicon, LLM inference is almost strictly memory bandwidth bound. The M3 Max might have 400 GB/s of bandwidth.

- **FP16 Model:** A 70B parameter model at 16-bit precision requires ~140GB of memory. Loading these weights once consumes significant bandwidth.
- **4-bit Quantization:** Reducing precision to 4-bits reduces the model size to ~35GB. This means the hardware can load the weights 4x faster (or load 4x more data in the same time).

**Insight:** For concurrent queries, quantization is not just about saving RAM; it is a throughput multiplier. By using 4-bit quantized models (which mlx-lm handles natively via Grouped quantization), you effectively quadruple the theoretical maximum token generation rate for a given batch size compared to FP16.[11]

## 4.2 mx.compile: The Fusion Reactor

Python overhead is real. Dispatching operations layer-by-layer, token-by-token in a loop creates "launch latency." mx.compile eliminates this by tracing the execution graph, fusing operations (e.g., combining MatMul + Bias + ReLU into a single GPU kernel), and reducing the communication overhead between Python and Metal.

**The Dynamic Shape Challenge:**

mx.compile is most effective when input shapes are static. If the batch size changes every iteration (e.g., from 3 concurrent users to 4, then back to 3), mx.compile may trigger a **recompilation**. Recompilation is expensive and causes stuttering.

**Optimization Strategy: Bucketing and Padding**

To optimize concurrent MLX queries, do not pass arbitrary batch sizes to the compiled function. Instead, use **Bucketing**.

- Define a set of fixed batch sizes: ``.
- If you have 3 active requests, pad the input tensor to size 4 using dummy data (masked out via the attention mask).
- Compile a separate graph for each bucket size. This ensures that the compiler always sees a known shape, keeping the inference path "warm" and optimized. Similarly, sequence lengths should be bucketed (e.g., pad to nearest 128 tokens) to avoid recompiling as the sequence grows.[3]

## 4.3 Speculative Decoding in a Batch

Speculative decoding involves using a small "draft" model to guess the next $N$ tokens, which are then verified by the large "target" model in a single parallel step.

- **Concurrency Synergy:** In a batch setting, speculative decoding can be powerful but complex. If the draft model is accurate, the batch jumps forward by multiple tokens.
- **MLX Implementation:** MLX supports speculative decoding natively. For concurrent queries, applying speculation per-sequence in the batch is feasible but requires careful management of the verification step to ensure that a rejection in one sequence doesn't stall the others.

# 5. Leveraging Batch Prompting and KV Cache

The user asks specifically about leveraging "batch prompting and cache." This is the mechanics of the engine.

## 5.1 Batch Prompting Mechanics

Batch prompting is the act of feeding multiple distinct prompts into the model simultaneously.

- **Tensor Layout:** If you have 3 prompts, you tokenize them into a list of lists.
  - Prompt A: `` (Length 3)
  - Prompt B: `` (Length 5)
- **Padding:** Tensors must be rectangular. You must pad Prompt A to length 5: [P, P, 101, 20, 5] (Left-padding is crucial for generation so the last token aligns).

- **Attention Mask:** You must construct a mask `` for Prompt A so the model ignores the padding.

**Impact:** Processing these together allows the GPU to perform Matrix-Matrix multiplications (GEMM) rather than Matrix-Vector (GEMV). GEMM is significantly more efficient on Apple's Neural Engine and GPU cores, leading to higher Token/Sec throughput per watt.[15]

## 5.2 The KV Cache: Anatomy and Optimization

The KV Cache stores the Key and Value matrices for every token generated so far. Without it, generating the 100th token would require re-processing the previous 99 tokens.

**Structure in MLX:**
The cache is typically a list of mx.array tuples (one per layer). For a batch of size $B$, sequence length $L$, and hidden dimension $D$, the cache size is $2 \times L \times B \times D$ (times layers).

**Handling "Persistent" Cache:**
In standard mlx-lm usage, the cache is created, used for one generation, and discarded. For concurrent systems (like a chatbot), we need **Persistent Caching**.

- **Session Management:** The server must maintain a dict mapping session_id to its specific KVCache object.
- **Context Swapping:** When a user sends a follow-up message, the server retrieves their existing cache. The new prompt tokens are appended to the computation, using the existing cache for context. This avoids re-computing the "history" of the chat.[12]

## 5.3 Prefix Caching: The "System Prompt" Optimization

Most concurrent queries in an agentic system share a common System Prompt (e.g., "You are a helpful assistant..."). Re-computing the KV cache for this prefix for every user is wasteful.

**Strategy:**
1. Compute the KV Cache for the System Prompt *once* at startup. Store it as a "Template Cache".
2. When a new request arrives, deep copy the Template Cache into the new request's session cache.
3. Start generation from the end of the system prompt. This reduces the "Time to First Token" (TTFT) for new requests significantly, as the heavy "prefill" of the system prompt is effectively skipped.[16]

## 5.4 Advanced Cache Management: PagedAttention

While mlx-lm uses contiguous memory arrays for caching (which can lead to fragmentation), the industry standard (pioneered by vLLM) is **PagedAttention**—breaking the cache into non-contiguous blocks.

- **In MLX:** As of early 2026, native PagedAttention is not the default in mlx-lm. However, community implementations (like vllm-mlx) have ported this logic. Implementing this manually requires managing a "Page Table" that maps logical token positions to physical memory blocks in an mx.array. This is advanced "Custom Handling" territory

but essential for maximizing memory efficiency in high-load scenarios.[17]

# 6. Gap Analysis: mlx-lm vs. Custom Engineering

To effectively plan a project, one must know what is given and what must be built.

## 6.1 What is already in mlx-lm?

- **Model Architectures:** Robust, optimized implementations of Llama 3, Mistral, Qwen, Phi, and MoE models. You do *not* need to write model code.
- **Weight Loading:** Seamless loading from Hugging Face, including handling safetensors and sharding.
- **Quantization:** High-performance kernels for 4-bit and 8-bit inference.
- **Basic Batching:** The generate() function accepts a list of prompts and handles padding/tokenization internally for a *static* batch.
- **Simple Server:** mlx_lm.server provides an OpenAI-compatible API, but it is generally a reference implementation suited for single-user or low-concurrency use cases (often serializing requests).[18]
- **Cache Classes:** KVCache and RotatingKVCache primitives are provided.[20]

## 6.2 What Must Be Handled "Ourselves" (Custom Engineering)

- **The Continuous Batching Scheduler:** The logic to dynamically insert/remove requests from a running batch loop is **not** in the core high-level API. You must build the while loop that manages the queue and calls model() directly.
- **Fairness Algorithms:** Algorithms like "FairBatching" (which balances prefill vs. decode priority to prevent starvation) must be implemented in the scheduler logic.[21]
- **Persistent Session Management:** The logic to map HTTP sessions to specific KV cache objects in memory and manage their lifecycle (eviction, LRU caching of sessions) is application-level logic.
- **Advanced Bucketing:** While mx.compile exists, the logic to bucket requests into fixed shapes to optimize compilation is a custom wrapper you must write around the model.
- **Structured Generation (Guided Decoding):** While tools like outlines exist, integrating strict JSON schema enforcement into a *concurrent* batch requires custom LogitsProcessor implementations that can handle a batch of requests where each request might have a *different* schema constraints.

# 7. Deep Dive: The generate_step Function

The atomic unit of customization is the generate_step function. Understanding this is mandatory for building custom schedulers.
**Source Code Analysis:**
Typically located in mlx_lm.utils, its signature resembles:

Python

```
def generate_step(prompt: mx.array, model: nn.Module, temp: float = 0.0,...):
    #...
```

It takes a tensor of input tokens and the model. It performs one forward pass.

**Usage in Concurrency:**
You should not use the stock generate_step for a complex server because it often assumes a single set of sampling parameters (temperature, top-p) for the whole batch. In a real server, User A might want temp=0.7 and User B temp=0.2.

**Refactored Logic:**
You must implement a custom step function that:
1. Takes the batch of inputs.
2. Runs the forward pass: logits, cache = model(inputs, cache).
3. **Splits** the logits.
4. Applies different sampling strategies to different rows of the logits tensor based on user preferences.
5. Returns the vector of next tokens.

This level of granularity is what distinguishes a "demo" server from a "production" server.

# 8. Case Study: Optimization with Structured Output

A rapidly growing requirement is **Structured Output** (forcing the LLM to output valid JSON).

- **Challenge:** Validating JSON requires a Finite State Machine (FSM) masking invalid tokens.
- **Concurrent Challenge:** If Request A needs JSON and Request B needs Python code, the batching engine needs to apply different masks to different batch indices.
- **Solution:** Integration with libraries like outlines. The scheduler must carry a "Logits Processor" state for each request. Before sampling, the logits for index i are passed through the processor for Request i. This is computationally cheap on the CPU but requires careful plumbing in the loop.[22]

# 9. Conclusion

Achieving concurrent queries with MLX in 2026 is a solved engineering problem, but it is not a "one-line import." It requires a comprehensive architectural approach that leverages Apple's Unified Memory for zero-copy data management and Python's asyncio for orchestrating the inference loop.

**Strategic Recommendations:**
1. **Adopt Continuous Batching:** Abandon static batching. Implement a dynamic scheduler that keeps the GPU saturated.
2. **Optimize the Cache:** Use Prefix Caching for system prompts and manage KV lifecycles explicitly.
3. **Compile with Intent:** Use mx.compile but protect it with input bucketing to prevent

recompilation thrashing.
4.  **Quantize:** Use 4-bit quantization as the default to maximize memory bandwidth utilization.

By following these principles, developers can unlock the full capability of Apple Silicon, transforming the Mac into a highly capable, local inference node that rivals dedicated cloud instances in efficiency and responsiveness.

## Data Tables and Comparisons

| Feature | mlx-lm Native Support | Custom Engineering Required |
|---|---|---|
| **Basic Inference** | Yes (generate, stream_generate) | No |
| **Static Batching** | Yes (via list input) | No |
| **Continuous Batching** | No | **Yes** (Scheduler Logic) |
| **Quantization** | Yes (4-bit, 8-bit) | No |
| **KV Cache** | Yes (Basic Objects) | **Yes** (Persistence/Session Mgmt) |
| **PagedAttention** | No (Core Lib) | **Yes** (or use vllm-mlx) |
| **Structured Output** | Basic | **Yes** (Concurrent FSM masking) |
| **Multi-GPU Pipeline** | Basic (mx.distributed) | **Yes** (Cluster Orchestration) |

*Table 1: Capability Matrix for Concurrency in MLX (Jan 2026)*

| Optimization | Mechanism | Impact on Concurrency |
|---|---|---|
| **Unified Memory** | Zero-copy CPU<->GPU access | **High**: Allows overlap of pre-processing and compute. |
| **4-bit Quantization** | Reduces model weight size | **Critical**: Quadruples available memory bandwidth for batching. |
| **Prefix Caching** | Reuses KV cache for shared prompts | **High**: Eliminates prefill latency for system prompts. |
| **Bucketing** | Pads inputs to fixed shapes | **High**: Prevents mx.compile recompilation stalls. |

*Table 2: Key Optimization Vectors*
1

## Works cited

1.  Exploring LLMs with MLX and the Neural Accelerators in the M5 GPU, accessed January 27, 2026, https://machinelearning.apple.com/research/exploring-llms-mlx-m5
2.  ml-explore/mlx: MLX: An array framework for Apple silicon - GitHub, accessed

January 27, 2026, https://github.com/ml-explore/mlx

3. Compilation — MLX 0.30.4 documentation, accessed January 27, 2026, https://ml-explore.github.io/mlx/build/html/usage/compile.html

4. python - multiprocessing vs multithreading vs asyncio - Stack Overflow, accessed January 27, 2026, https://stackoverflow.com/questions/27435284/multiprocessing-vs-multithreading-vs-asyncio

5. Deep Dive into Multithreading, Multiprocessing, and Asyncio | by Clara Chong - Medium, accessed January 27, 2026, https://medium.com/data-science/deep-dive-into-multithreading-multiprocessing-and-asyncio-94fdbe0c91f0

6. Question about running `mx.eval` in separate threads · ml-explore mlx · Discussion #1448 - GitHub, accessed January 27, 2026, https://github.com/ml-explore/mlx/discussions/1448

7. Thread safety: Ongoing issue for thread safety in MLX · Issue #2133 ..., accessed January 27, 2026, https://github.com/ml-explore/mlx/issues/2133

8. MLX_LM Continuous Batching. Apple Silicon is one of the easiest... | by Melkor - Medium, accessed January 27, 2026, https://medium.com/@clnaveen/mlx-lm-continuous-batching-e060c73e7d98

9. mlx-openai-server - PyPI, accessed January 27, 2026, https://pypi.org/project/mlx-openai-server/

10. How continuous batching enables 23x throughput in LLM inference while reducing p50 latency - Anyscale, accessed January 27, 2026, https://www.anyscale.com/blog/continuous-batching-llm-inference

11. WWDC 2025 - Explore LLM on Apple silicon with MLX - DEV Community, accessed January 27, 2026, https://dev.to/arshtechpro/wwdc-2025-explore-llm-on-apple-silicon-with-mlx-1if7

12. MLX batch generation is pretty cool! : r/LocalLLaMA - Reddit, accessed January 27, 2026, https://www.reddit.com/r/LocalLLaMA/comments/1fodyal/mlx_batch_generation_is_pretty_cool/

13. PyTorch and MLX for Apple Silicon | Towards Data Science, accessed January 27, 2026, https://towardsdatascience.com/pytorch-and-mlx-for-apple-silicon-4f35b9f60e39/

14. Writing Fast MLX - GitHub Gist, accessed January 27, 2026, https://gist.github.com/awni/4beb1f7dfefc6f9426f3a7deee74af50

15. Batched KV caching in MLX for fast parallel generation (1000+ tok/s throughput for Gemma-2B + Phi-3-mini on M3 Max) - Reddit, accessed January 27, 2026, https://www.reddit.com/r/LocalLLaMA/comments/1dn2b4n/batched_kv_caching_in_mlx_for_fast_parallel/

16. Plans for batch KV cache? · Issue #548 · ml-explore/mlx-lm - GitHub, accessed January 27, 2026, https://github.com/ml-explore/mlx-lm/issues/548

17. vLLM-MLX: Native Apple Silicon LLM inference - 464 tok/s on M4 Max :

r/LocalLLaMA, accessed January 27, 2026,
https://www.reddit.com/r/LocalLLaMA/comments/1qeley8/vllmmlx_native_apple_silicon_llm_inference_464/

18. mlx-textgen - PyPI, accessed January 27, 2026,
https://pypi.org/project/mlx-textgen/0.1.4/

19. ml-explore/mlx-lm: Run LLMs with MLX - GitHub, accessed January 27, 2026,
https://github.com/ml-explore/mlx-lm

20. Releases · ml-explore/mlx-lm - GitHub, accessed January 27, 2026,
https://github.com/ml-explore/mlx-lm/releases

21. FairBatching: Fairness-Aware Batch Formation for LLM Inference - arXiv, accessed
January 27, 2026, https://arxiv.org/html/2510.14392v1

22. MLX batched/continous inference with structured outputs : r/LocalLLaMA -
Reddit, accessed January 27, 2026,
https://www.reddit.com/r/LocalLLaMA/comments/1qjpjt6/mlx_batchedcontinous_inference_with_structured/

23. LM Studio 0.3.4 ships with Apple MLX, accessed January 27, 2026,
https://lmstudio.ai/blog/lmstudio-v0.3.4