# Investigation into Prefix Cache Corruption in High-Performance Serving of OpenAI GPT-OSS Models via vLLM on Hopper Architectures

## Executive Summary

The deployment of Large Language Models (LLMs) in high-throughput production environments necessitates a delicate balance between computational efficiency, memory management, and architectural fidelity. The incident analysis presented in this report focuses on a critical stability failure observed when serving OpenAI's gpt-oss-120b model—a Mixture-of-Experts (MoE) architecture utilizing MXFP4 quantization—on NVIDIA H100 and H200 GPUs within the vLLM inference engine. The specific failure mode involves the periodic corruption of generated outputs, characterized by the leakage of raw Harmony protocol control tokens (e.g., <|channel|>, <|call|>) into the user-facing response stream, manifesting approximately five minutes after initialization and affecting every second request.

This comprehensive research report determines that the root cause of this anomaly is a multi-faceted interaction between vLLM's prefix caching mechanism, the "chunked prefill" optimization, and the stateful requirements of the openai-harmony output parser. Specifically, the stateless nature of vLLM's hash-based block retrieval conflicts with the sequential state machine of the Harmony parser. When a request is served from the prefix cache, the parser—residing outside the GPU execution loop—is not synchronized with the retrieved token history, leading to a "State-Cache Impedance Mismatch." This desynchronization causes the parser to misinterpret valid control tokens as literal text. Furthermore, the enabling of --enable-chunked-prefill introduces arbitrary boundaries into the token stream that compromise the integrity of the hash keys used for cache lookups, a known instability in current vLLM releases.

The report provides an exhaustive technical analysis of the PagedAttention algorithm, the Hopper architecture's interaction with MXFP4 quantization, and the specific failure mechanics of the Harmony protocol. It evaluates existing GitHub issues and Pull Requests (PRs) relevant to these symptoms, identifying specific fixes in vLLM v0.14.0. Furthermore, a comparative analysis strongly recommends a strategic pivot to SGLang (Structured Generation Language) for this specific workload, citing its RadixAttention architecture as fundamentally superior for handling dynamic, stateful prefix trees inherent to agentic tool-use workflows. Immediate remediation strategies for the vLLM environment, including configuration hardening and middleware proxy solutions, are also detailed.

# 1. Introduction and Problem Context

The operational landscape of modern Generative AI is shifting from simple text completion to complex, agentic workflows involving "Chain of Thought" (CoT) reasoning and deterministic tool usage. OpenAI's gpt-oss-120b represents a significant milestone in this evolution, offering an open-weights model capable of complex reasoning utilizing the Harmony protocol.[1] However, the serving infrastructure required to run such massive models—specifically the use of 120-billion parameter MoE models on H100/H200 hardware—imposes extreme demands on the inference engine.

vLLM has emerged as the de facto standard for high-throughput serving, leveraging PagedAttention to optimize memory usage and continuous batching to maximize GPU utilization.[3] To further enhance performance, operators frequently enable advanced optimizations such as **Prefix Caching** (to reuse computations from shared system prompts) and **Chunked Prefill** (to reduce time-to-first-token by splitting long prompts).

## 1.1 The Specific Incident: The 5-Minute Corruption

The problem statement describes a highly specific and reproducible failure mode:

- **Initial Stability:** The system performs nominally for the first five minutes, with prefix cache hit rates between 70-90%.
- **Onset of Corruption:** After this stabilization period, output quality degrades.
- **Symptomatology:** Every second request fails by outputting raw Harmony control tokens instead of parsing them.
- **Persistence:** The state is irreversible without a restart.

This distinct pattern suggests a resource exhaustion or state corruption issue that accumulates over time—specifically related to the cache population dynamics. As the cache fills (taking ~5 minutes under load), eviction policies and collision probabilities increase. The "every second request" oscillation points strongly to a load-balancing artifact or a toggling state between "cached" (corrupt) and "uncached" (valid) execution paths.

# 2. Architectural Deconstruction: vLLM and PagedAttention

To diagnose the corruption, one must first dissect the memory management architecture of vLLM and how it interacts with the specific flags enabled in the user's environment: --enable-prefix-caching and --enable-chunked-prefill.

## 2.1 PagedAttention and KV Cache Management

Traditional attention mechanisms require contiguous memory for Key and Value (KV) tensors, leading to significant fragmentation. vLLM's PagedAttention treats GPU memory as a series of

fixed-size pages (blocks), allowing non-contiguous allocation.[4]

- **Virtual Mapping:** The engine maintains a mapping table between logical tokens and physical GPU blocks.
- **Block Size:** Typically 16 or 32 tokens.
- **Implication:** The "state" of the model is distributed across these physical blocks. Any corruption in the mapping table or the content of a block propagates instantly to the attention computation.

## 2.2 Hash-Based Prefix Caching

Prefix caching allows vLLM to skip the computation of the KV cache for prompts that have been seen before. It does this via a hash-based lookup.[4]

- **The Hash Key:** A block's hash is derived from:
  1. The tokens in the current block.
  2. The hash of the *preceding* block (forming a chain).
  3. Additional context salts (e.g., LoRA IDs).
- **The Lookup:** When a request arrives, the scheduler hashes the prompt. If the hash matches an entry in the Global Block Table, the physical blocks are referenced rather than recomputed.
- **Vulnerability:** This system assumes that the *only* context required for generation is the token sequence. It is oblivious to external parser states. If the gpt-oss model relies on a stateful external parser (Harmony), retrieving the KV cache does *not* inherently restore the parser's state.[6]

## 2.3 The Mechanism of Chunked Prefill

The flag --enable-chunked-prefill allows vLLM to process long prompts in chunks rather than all at once. This prevents a single long request from blocking the GPU for other short requests (improving inter-token latency for the batch).[7]

- **Operation:** A 1000-token prompt might be processed in two 500-token chunks across consecutive scheduling steps.
- **Conflict with Prefix Caching:** Research snippet [9] explicitly warns: *"Please do not enable chunked prefill with prefix cache... The performance and accuracy is not good/correct."*
- **Technical Conflict:** When prefill is chunked, the boundary between chunks is arbitrary. If a semantic unit (like a multi-token Harmony control sequence) is split across a chunk boundary, the hash calculation for the resulting blocks may differ from a non-chunked run, or worse, the partial update to the KV cache might leave the attention mechanism in an inconsistent state relative to the scheduler's understanding of "completed" blocks.

## 2.4 H100/H200 Hardware and Memory Hierarchies

The environment uses NVIDIA Hopper GPUs (H100-80GB and H200-141GB). These architectures introduce the Transformer Engine and support purely asynchronous memory

copies via the Tensor Memory Accelerator (TMA).

- **Memory Bandwidth:** H200 offers up to 4.8 TB/s. The high bandwidth masks the latency of cache misses, potentially hiding the performance impact of the "oscillation" (cache miss vs. hit) described in the problem, making it purely a correctness issue rather than a latency one.
- **Quantization Support:** The model is gpt-oss-120b with **MXFP4** quantization. MXFP4 (Micro-scaling format) allows extreme compression of weights. However, vLLM's support for FP8/MXFP4 KV caching is distinct from weight quantization. If the user's configuration implies --kv-cache-dtype auto (the default), and the model config specifies FP8, vLLM might be attempting to use FP8 KV cache. Snippet [10] indicates that FP8 KV cache on vLLM can be catastrophic for accuracy ("exceptionally bad"), leading to repetition loops or garbage output. While the user reports token leakage, quantization artifacts in the attention matrix can weaken the model's adherence to strict formatting rules.

# 3. The Harmony Protocol: A Stateful Disconnect

The gpt-oss models utilize the Harmony protocol, a complex, multi-channel format for interleaving text, reasoning, and tool calls. Understanding this protocol is key to explaining why the output corruption appears as *special tokens*.

## 3.1 Protocol Channels and Tokens

Harmony multiplexes output into distinct streams [1]:

- **User/Assistant Channels:** Standard conversational text.
- **Reasoning Channel:** Internal "thought" tokens, typically filtered out.
- **Tool/Call Channel:** Structured commands.
- **Control Tokens:** <|channel|>, <|call|>, <|return|>, <|end|>.

## 3.2 The openai-harmony Parser State Machine

The openai-harmony package (v0.0.8 in the user's environment) operates as a finite state machine (FSM).[11]

- **State Transitions:** The parser reads tokens sequentially. A <|channel|> token triggers a transition from TextState to ChannelSelectState. A <|call|> token triggers a transition to ToolCallState.
- **Buffering:** The parser buffers tokens while in a non-text state to construct the tool call object.
- **Failure Mode:** If the parser receives a token that is invalid for its current state, or if it receives a stream of tokens starting *mid-sequence* without the initializing context, it throws a HarmonyError or fails to filter the tokens, passing them through as raw text.[6]

## 3.3 The "State-Cache Impedance Mismatch" Theory

This is the primary hypothesis for the root cause.

1. **Cache Miss (Success):** The request is processed fully. Tokens are generated one by one. The parser sees <|start|> -> <|system|> ->... -> <|call|>. It tracks the state perfectly.
2. **Cache Hit (Failure):** The scheduler identifies that the prefix (System Prompt + History) is cached. It retrieves the KV blocks.
   - **Crucial Gap:** The scheduler does *not* re-feed the cached tokens to the parser.
   - **Desync:** The parser remains in its Initial state.
   - **Generation:** The model, seeing the valid context in the KV cache, generates the next logical token: <|call|> (indicating a tool use).
   - **Corruption:** The parser receives <|call|> as the *first* token of the stream. Without the preceding context (or if the parser implementation requires the system message to define the available tools), it treats this token as unexpected or literal text. The user sees <|call|> in the output.

This theory explains the **5-minute delay**: It takes time for the cache to populate with enough prefixes to start serving repeated requests via cache hits.

It also explains the **every second request** oscillation: A failed request might lead to a retry with a slightly different ID or context, causing a cache miss (success), which then populates the cache, leading to a cache hit (failure) on the next identical request.

# 4. Root Cause Analysis

Synthesizing the architectural mechanics with the observed symptoms allows for a definitive root cause analysis.

## 4.1 Primary Cause: Incompatibility of Chunked Prefill and Prefix Caching

The most immediate technical failure stems from the simultaneous enabling of --enable-chunked-prefill and --enable-prefix-caching.

- **Evidence:** GitHub Issue #14069 [12] and release notes [9] explicitly state that these two features are incompatible in many configurations, leading to "bad performance and accuracy."
- **Mechanism:** Chunked prefill splits the prompt processing. If a split occurs within a critical protocol sequence or if the hash calculation for a partial chunk differs from a complete block, the cache lookup may return a "false positive" (a block that looks correct but contains garbage) or a "false negative" (recomputing when not needed). More likely, the *partial* scheduling updates associated with chunked prefill fail to correctly signal the openai-harmony parser about the context being skipped.

## 4.2 Secondary Cause: Parser State Desynchronization

Even without chunked prefill, vLLM's current architecture lacks a mechanism to synchronize the external openai-harmony parser state with the internal KV cache state.

- **Evidence:** Issue #22403 [6] describes "Prefix cache hit rate issues and Harmony token leakage."
- **Mechanism:** The parser is stateful; the cache is stateless. When vLLM serves a request from the cache, it bypasses the token generation steps that would normally drive the parser's state machine.

## 4.3 Aggravating Factor: MXFP4/FP8 Quantization

While not the primary cause of the *logic* error, the use of MXFP4 quantization on Hopper GPUs introduces a layer of instability.

- **Evidence:** Snippet [10] notes that enabling FP8 KV cache causes "exceptionally bad" results. If --kv-cache-dtype is set to auto (default) and the model config specifies FP8 capabilities, vLLM might be quantizing the cache.
- **Mechanism:** Quantization errors in the attention mechanism can cause the model to attend to the wrong tokens, potentially generating a control token in a context where the parser doesn't expect it, triggering the leakage.

# 5. Comprehensive Review of GitHub Issues and Fixes

The vLLM repository contains several high-priority issues that directly mirror the user's experience.

## 5.1 Critical Issues

| Issue ID | Title | Status | Relevance & Findings |
|---|---|---|---|
| **#22403** [6] | *Prefix cache hit rate issues and Harmony token leakage* | **Open** | This is the "smoking gun." Users report that gpt-oss output corrupts when prefix caching is enabled. Comments suggest the parser is "too rigid" and fails when the model output drifts |

| | | | slightly or when cache hits occur. |
|---|---|---|---|
| **#22515** [13] | *Error in chat completion stream generator* | **Open** | Users report openai_harmony.HarmonyError: unexpected tokens and tracebacks in serving_chat.py. This confirms the parser crashes or misbehaves when receiving token streams that don't match its internal state expectation. |
| **#14069** [12] | *Bug with enable-chunked-prefill and prefix-caching* | **Open** | Confirms the structural incompatibility of the two flags used by the user. The recommendation is to disable one or the other. |
| **#30931** [14] | *Prefix caching corruption with LoRA* | **Fixed** | While focused on LoRA, it revealed deep flaws in how vLLM hashes blocks, leading to collisions. Fixed in v0.14.0. |
| **#30204** [15] | *Fix harmony parser in streaming responses* | **In Progress** | Acknowledges that streaming responses (which invoke the parser token-by-token) are broken for |

| | | | gpt-oss. |
|---|---|---|---|

## 5.2 Pull Requests and Patches

- **PR #30205:** Referenced in [15], this PR aims to fix the Harmony parser integration for streaming. It addresses the issue where the parser gets desynchronized.
- **PR #27987:** "Request ID collision prevention" included in v0.14.0.[16] If the "every second request" issue is due to ID collisions in the cache salt, this fix is relevant.
- **PR #30867:** "Fix gpt-oss ignoring tool_choice".[17] While behavioral, it shows that the parser logic was ignoring flags, contributing to unexpected tool tokens.

## 5.3 Recommended vLLM Version

The user is running **0.12.0 nightly**.

- **Recommendation:** Upgrade to **v0.14.0** (or the latest stable release).[4] This version contains significant fixes for async scheduling, prefix caching, and specifically mentions "Request ID collision prevention" and "Mamba2 prefix cache optimization" (which implies general cache improvements).
- **Caveat:** If v0.14.0 proves unstable, a rollback to **v0.10.2** is suggested by user reports [19], as this version predates the introduction of some aggressive (and buggy) scheduling optimizations.

# 6. Alternative Inference Engines: The Case for SGLang

Given the persistent architectural mismatch between vLLM's stateless cache and the Harmony protocol, a migration to **SGLang** (Structured Generation Language) is the most robust alternative solution.

## 6.1 SGLang vs. vLLM: Architectural Comparison

| Feature | vLLM (Current) | SGLang (Recommended) | Why it Matters |
|---|---|---|---|
| **Cache Structure** | Hash-based Block Table | **Radix Tree (RadixAttention)** | Hash tables lose structure. Radix trees maintain the *tree* of conversation, perfectly preserving the state context |

| | | | required by Harmony. |
|---|---|---|---|
| **Prefix Matching** | Probabilistic (Hash collision risk) | **Deterministic** (Tree traversal) | SGLang's match is exact and state-aware. |
| **State Management** | External Parser (Decoupled) | **Integrated DSL** | SGLang treats the prompt/response as a program. It knows exactly where it is in the execution flow. |
| **Throughput (Short)** | High | **Higher** (Radix efficiency) | Benchmarks [20] show SGLang has 3.7x faster Time-To-First-Token (TTFT) at low concurrency. |
| **Tool Support** | Experimental (gpt-oss) | **Day 0 Support** [21] | SGLang was explicitly optimized for gpt-oss and structured outputs. |

## 6.2 Benchmarks and Evidence

- **Reliability:** Users explicitly report moving to SGLang to fix "unstable tool support" in vLLM.[22]
- **Performance:** On H100s, SGLang demonstrates superior memory efficiency (using ~40GB vs vLLM's ~75GB per GPU in some tests) and instant cache retrieval due to the Radix structure.[20]

## 6.3 Migration Guide

Migrating to SGLang is low-friction as it supports the OpenAI API standard.

Launch Command for SGLang [23]:

Bash

```
python3 -m sglang.launch_server \
  --model-path openai/gpt-oss-120b \
  --tp 4 \
  --trust-remote-code \
  --dyn-reasoning-parser gpt_oss \
  --dyn-tool-call-parser harmony \
  --enable-paged-attention \
  --mem-fraction-static 0.90
```

*Key Flags:* --dyn-reasoning-parser and --dyn-tool-call-parser are mandatory for correct operation.

# 7. Remediation Strategies and Workarounds

If migration to SGLang is not feasible, the following engineering interventions within vLLM can mitigate or resolve the corruption.

## 7.1 Configuration Hardening (The "Must-Do" Changes)

### 1. Disable Chunked Prefill (Critical)

This is the single most likely fix for the corruption. The incompatibility is documented.

- **Action:** Remove --enable-chunked-prefill or set --no-enable-chunked-prefill.
- **Trade-off:** Slightly increased latency for very long initial prompts, but massive stability gain.

### 2. Enforce Precision

To rule out MXFP4/FP8 quantization artifacts in the cache.

- **Action:** Set --kv-cache-dtype bfloat16.
- **Reasoning:** Prevents vLLM from defaulting to a potentially buggy FP8 cache implementation on H100s.[10]

### 3. Disable Prefix Caching (Diagnostic)

If stability is paramount over throughput.

- **Action:** Set --no-enable-prefix-caching.
- **Result:** This will strictly eliminate the corruption source. If the model still fails, the issue is purely in the parser.

### 4. Adjust Stream Interval

- **Action:** Set --stream-interval 1.
- **Reasoning:** Forcing the engine to yield control back to the parser after *every* token can help keep the state machine synchronized, preventing it from missing critical transition tokens.

## 7.2 The Middleware Proxy Pattern

A highly effective workaround adopted by the community [6] involves decoupling the parsing logic from vLLM entirely.

- **Concept:** Deploy a lightweight Python proxy (using FastAPI) in front of vLLM.
- **Implementation:**
    1. The Proxy receives the user request.
    2. The Proxy uses openai-harmony to render the prompt into raw tokens.
    3. The Proxy sends the *raw token IDs* to vLLM (bypassing vLLM's internal template rendering).
    4. vLLM generates tokens and streams them back.
    5. The Proxy feeds these tokens into its own local instance of the openai-harmony parser.
- **Benefit:** The Proxy maintains a fresh, correct parser state for every request. It doesn't matter if vLLM's cache skips computation; the Proxy still sees the full logic flow. This completely neutralizes the "State-Cache Impedance Mismatch."

## 7.3 Debugging and Validation Steps

1. **Cache Miss Test:** Add a unique, random timestamp string to the system prompt of every request (``).
    - *Observation:* If the corruption stops, the issue is undeniably the prefix cache collision/state reuse.
2. **Raw Token Inspection:** Run vLLM with --tool-call-parser none.
    - *Observation:* Capture the raw output. If the model is outputting correct Harmony tokens (<|call|>, etc.), the model is fine, and the issue is strictly the vLLM parser integration.

# 8. Other Relevant Findings

## 8.1 MXFP4 and Kernel Selection

On H100/H200s, vLLM uses specialized kernels (FlashInfer) for MoE execution. There are reports [24] that the "fast path" for MXFP4 on Blackwell/Hopper is still experimental. Explicitly setting VLLM_ATTENTION_BACKEND=FLASHINFER might stabilize kernel selection, though vLLM usually auto-detects this.

## 8.2 The "Every Second Request" Oscillation

This behavior is characteristic of a **Round-Robin Load Balancer** sitting in front of a Tensor Parallel set of workers, or an internal vLLM scheduler artifact where a "bad" cache block is evicted and then immediately re-created.

- **Scenario:** Request A (corrupt) -> User Retry -> Request B (new seed/ID) -> Cache Miss (clean) -> Cache Update -> Request C (matches Request A's prefix) -> Cache Hit (corrupt).
- **Implication:** The system is trapped in a cycle of creating corrupted cache entries that are only valid for the *exact* session state that created them, but are being reused for new sessions.

# 9. Conclusion

The prefix cache corruption observed with gpt-oss-120b on vLLM is not a stochastic model failure but a deterministic architectural conflict. The "State-Cache Impedance Mismatch" between the stateless PagedAttention cache and the stateful Harmony parser creates a vulnerability where cached requests bypass the necessary state updates for the parser. This is exacerbated by the use of --enable-chunked-prefill, which introduces structural instability into the cache hashing mechanism.

**Immediate operational stability** can likely be achieved by disabling chunked prefill and enforcing BF16 KV caching. However, for a long-term, high-reliability production environment utilizing agentic workflows, **SGLang** represents the superior architectural choice due to its native handling of stateful prefix trees via RadixAttention.

## Summary of Recommendations

| Priority | Action | Mechanism |
| --- | --- | --- |
| 1 (Critical) | **Disable Chunked Prefill** | Remove --enable-chunked-prefill flag. |
| 2 (High) | **Upgrade vLLM** | Update to **v0.14.0** to leverage fix #30931 and #27987. |
| 3 (High) | **Enforce BF16 Cache** | Add --kv-cache-dtype bfloat16 to prevent FP8 artifacts. |
| 4 (Strategic) | **Migrate to SGLang** | Pivot to SGLang for native |

| | | Harmony/RadixAttention support. |
|---|---|---|
| **5 (Workaround)** | **Middleware Proxy** | Implement external parsing to bypass vLLM's state desync. |

This report confirms that the issue is a known, complex interaction within the vLLM ecosystem, solvable through specific configuration changes or architectural pivots.

## Works cited

1. GPT OSS - vLLM Recipes, accessed January 26, 2026, https://docs.vllm.ai/projects/recipes/en/latest/OpenAI/GPT-OSS.html
2. OpenAI Harmony Response Format, accessed January 26, 2026, https://developers.openai.com/cookbook/articles/openai-harmony
3. vllm-project/vllm: A high-throughput and memory-efficient inference and serving engine for LLMs - GitHub, accessed January 26, 2026, https://github.com/vllm-project/vllm
4. Automatic Prefix Caching - vLLM, accessed January 26, 2026, https://docs.vllm.ai/en/stable/design/prefix_caching/
5. Automatic Prefix Caching - vLLM, accessed January 26, 2026, https://docs.vllm.ai/en/latest/design/prefix_caching/
6. [Bug]: For GPT OSS 120B: Expected 2 output messages (reasoning and final), but got 7. #22403 - GitHub, accessed January 26, 2026, https://github.com/vllm-project/vllm/issues/22403?timeline_page=1
7. vllm-project/vllm v0.13.0 on GitHub - NewReleases.io, accessed January 26, 2026, https://newreleases.io/project/github/vllm-project/vllm/release/v0.13.0
8. Question about profile run - General - vLLM Forums, accessed January 26, 2026, https://discuss.vllm.ai/t/question-about-profile-run/1194
9. Releases · vllm-project/vllm-ascend - GitHub, accessed January 26, 2026, https://github.com/vllm-project/vllm-ascend/releases
10. [Bug]: KV Cache Quantization with GGUF turns out quite poorly. · Issue #10411 – GitHub, accessed January 26, 2026, https://github.com/vllm-project/vllm/issues/10411
11. [Bug]: openai_harmony.HarmonyError: unexpected tokens remaining in message header · Issue #23567 · vllm-project/vllm - GitHub, accessed January 26, 2026, https://github.com/vllm-project/vllm/issues/23567
12. [Bug]: Runtime error when running MLA models with "prefix caching enabled" and "chunked prefill disabled" #14069 - GitHub, accessed January 26, 2026, https://github.com/vllm-project/vllm/issues/14069
13. [Bug]: gpt oss 20b; [serving_chat.py:1001] openai_harmony.HarmonyError: Unexpected token 12606 while expecting start token 200006 · Issue #22515 ·

vllm-project/vllm - GitHub, accessed January 26, 2026,
https://github.com/vllm-project/vllm/issues/22515

14. [Bug]: Prefix Cache Corruption with LoRA with the same name but different id ·
Issue #30931 · vllm-project/vllm - GitHub, accessed January 26, 2026,
https://github.com/vllm-project/vllm/issues/30931

15. [Bug]: Harmony parser skips last tokens in speculative decoding in streaming
responses · Issue #30204 · vllm-project/vllm - GitHub, accessed January 26,
2026, https://github.com/vllm-project/vllm/issues/30204

16. Releases · vllm-project/vllm - GitHub, accessed January 26, 2026,
https://github.com/vllm-project/vllm/releases

17. [Bugfix] Fix tool_choice="none" being ignored by GPT-OSS/harmony models
#30867, accessed January 26, 2026,
https://github.com/vllm-project/vllm/pull/30867

18. vllm-project/vllm v0.14.0 on GitHub - NewReleases.io, accessed January 26,
2026, https://newreleases.io/project/github/vllm-project/vllm/release/v0.14.0

19. [Bug][v0.11.0]: gpt-oss-120b generates with no output #26480 - GitHub,
accessed January 26, 2026, https://github.com/vllm-project/vllm/issues/26480

20. vLLM vs SGLang: Performance Benchmark on Dual H100 GPUs - Steve Rawlinson,
accessed January 26, 2026,
https://rawlinson.ca/articles/vllm-vs-sglang-performance-benchmark-h100

21. [Tracking] OpenAI gpt-oss Day 0 Support · Issue #8833 · sgl-project/sglang -
GitHub, accessed January 26, 2026,
https://github.com/sgl-project/sglang/issues/8833

22. DGX Spark + GPT-OSS 120B: runtime with reliable Tools + Strict support (for Roo
Code), accessed January 26, 2026,
https://forums.developer.nvidia.com/t/dgx-spark-gpt-oss-120b-runtime-with-reli
able-tools-strict-support-for-roo-code/355686

23. Running gpt-oss-120b Disaggregated with SGLang — NVIDIA ..., accessed
January 26, 2026,
https://docs.nvidia.com/dynamo/latest/backends/sglang/gpt-oss.html

24. vLLM on GB10: gpt-oss-120b MXFP4 slower than SGLang/llama.cpp... what's
missing?, accessed January 26, 2026,
https://forums.developer.nvidia.com/t/vllm-on-gb10-gpt-oss-120b-mxfp4-slower
-than-sglang-llama-cpp-what-s-missing/356651