



Millstone AI Solutions

On-Prem AI Solutions

BENCHMARK REPORT

GLM-4.7-Flash

Performance Analysis on 1x H200 SXM

MODEL

Organization **Z.ai**
Parameters **30B**
Precision **BF16**

TEST HARDWARE

GPU **1x H200 SXM**
VRAM **141GB**
Engine **vLLM**

HIGHLIGHTS

458.3

Tok/s Peak Throughput
@ 4 Concurrent Requests

100.0%

Success Rate
Across All Scenarios

12+

Concurrent Users
@ 32K Context

Table of Contents

Executive Summary	3
Use Case Guidance	4
Performance Analysis	5
System Throughput	5
Per-User Generation Speed	6
Time to First Token	7
Capacity Analysis	8
Code Completion (1K Context)	8
Short-form Chatbot (8K Context)	9
General Chatbot (32K Context)	10
Long Document Processing (64K Context)	11
Automated Coding Assistant (96K Context)	12
Technical Deep Dive	13
Queue Wait Times	13
Per-User Prefill Speed	14
Inter-Token Latency	15
Scaling Efficiency	15
Power & Efficiency	16

Interactive Data Available Online

This report provides a static snapshot of benchmark results. For interactive charts with hover tooltips, exact data point values, and interpolated metrics, visit the full benchmark page:

MillstoneAI.com/inference-benchmark/glm-4.7-flash-bf16-1x-h200-sxm

Executive Summary

Infrastructure decisions require real performance data. This report measures user-facing performance, showing how many concurrent users a configuration can support at a given context length before performance degrades.

This benchmark evaluates **GLM-4.7-Flash** (Z.ai, 30B parameters, Mixture-of-Experts) running in BF16 precision on 1x H200 SXM (141GB VRAM).

Test parameters: Context lengths from 1K - 200K tokens. Concurrency from 1 - 4 requests. 1024 output tokens per request. No prompt caching. No speculative decoding. Full-precision KV cache.

[Benchmark methodology](#) →

Key Findings

Peak System Throughput	458.3 tok/s @ 4 concurrent requests, 1K context
TTFT Single Request	66ms (1K context) → 35.9s (200K context)
Generation Speed Single Request	196.7 tok/s (1K context) → 116.3 tok/s (200K context)
Chatbot Capacity	12+ concurrent requests @ 32K context
Throughput Scaling	2.9× from 1 to 4 concurrent requests
Success Rate	100.0% across 4.7K requests

Throughout this report, "**concurrent requests**" refers to simultaneous active requests. For applications with natural user pauses (chat interfaces, coding assistants), each request slot typically serves 4–5 active users.

RECOMMENDATIONS

Use Case Guidance

The table below maps this configuration's performance to common deployment scenarios. Capacity limits are where TTFT or generation speed falls below accepted thresholds for a comfortable user experience. Detailed charts and analysis for each use case are available in the Capacity Analysis section.

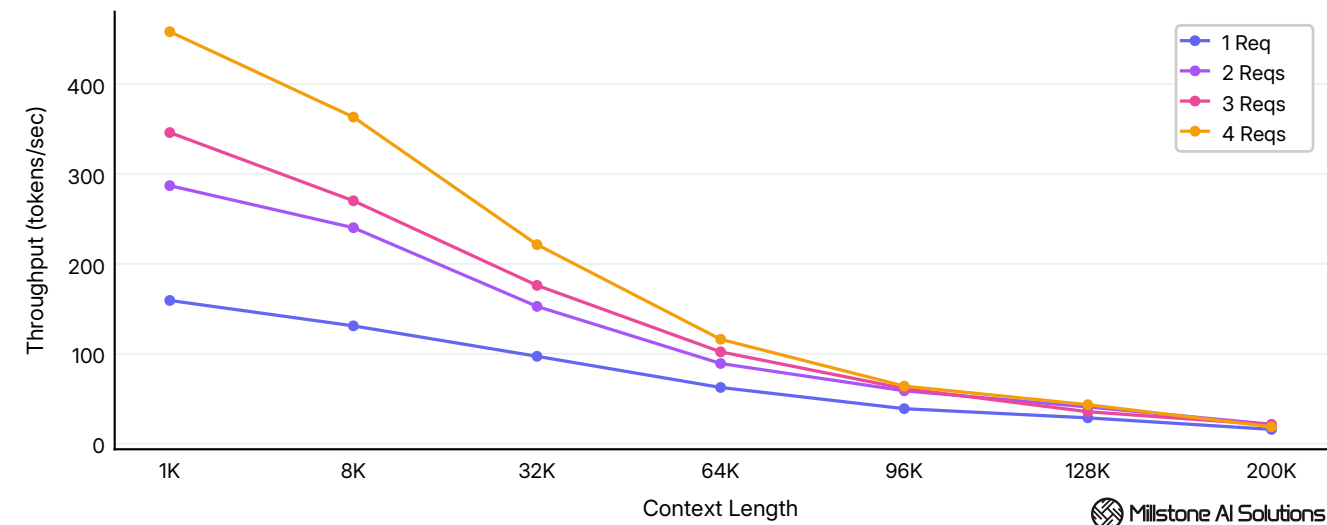
USE CASE	TTFT THRESHOLD	SPEED THRESHOLD	ANALYSIS
Code Completion	2s e2e	N/A	Supports 24 concurrent requests within accepted thresholds.
Short-form Chatbot	10s	10 tok/s	Supports 125+ concurrent requests with fast responses. Additional capacity likely available.
General Chatbot	8s	15 tok/s	Supports 12+ concurrent requests with fluid conversations. Significant additional capacity likely available.
Long Document Processing	12s	15 tok/s	Supports 4 concurrent requests within accepted thresholds.
Automated Coding Assistant	12s	20 tok/s	Best suited for single-user agentic workflows. For team environments, enable prompt caching or consider a smaller model.

The limits shown are conservative. Beyond these points, the system continues functioning with slower response times that may still be acceptable for your specific use case.

Want to validate your specific configuration? [Request a Custom Benchmark](#) →

System Throughput

Aggregate token generation across all concurrent requests. Measures output tokens only. Prompt tokens processed during prefill are excluded.



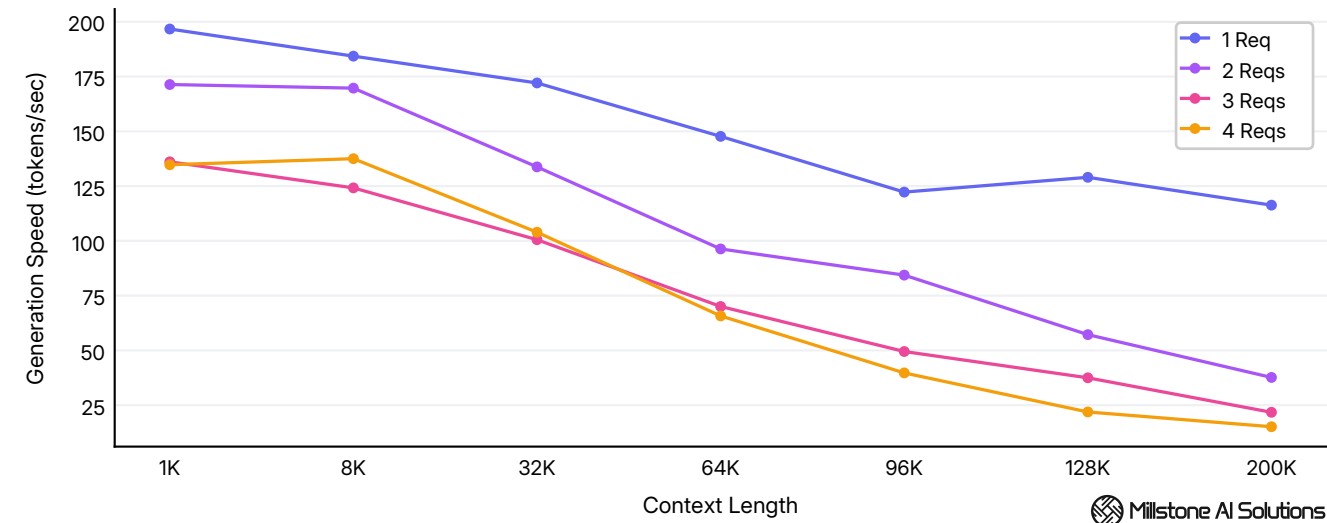
Average system throughput across 1K - 200K tokens context lengths at 1 - 4 concurrency levels.

CONDITION	THROUGHPUT
Peak (1K context, 4 requests)	458.3 tok/s
32K context, 4 requests	221.4 tok/s
200K context, 4 requests	18.8 tok/s

At peak throughput, this configuration produces approximately **1.6 million** tokens per hour. This is relevant for batch workloads like document processing, synthetic data generation, or offline analysis. Higher concurrency or shorter contexts can increase this further.

Per-User Generation Speed

Token generation rate experienced by each individual user. This is the speed at which text streams into their response, also referred to as "decode speed" or "decode throughput." As concurrency increases, per-user speed decreases since GPU resources are shared across requests.



Average per-user generation speed across 1K - 200K tokens context lengths at 1 - 4 concurrency levels.

How Fast is This?

SPEED	USER EXPERIENCE
< 15 tok/s	Slow; may be slower than reading speed
15–25 tok/s	Acceptable; keeps pace with reading
25–50 tok/s	Fast; exceeds reading speed
> 50 tok/s	Very fast; text appears nearly instantly

At **15.2 tok/s** (the lowest measured point: 200K context, 4 concurrent requests), this configuration stays at acceptable speeds across all tested scenarios. Single-user performance at 1K context reaches **196.7 tok/s**.

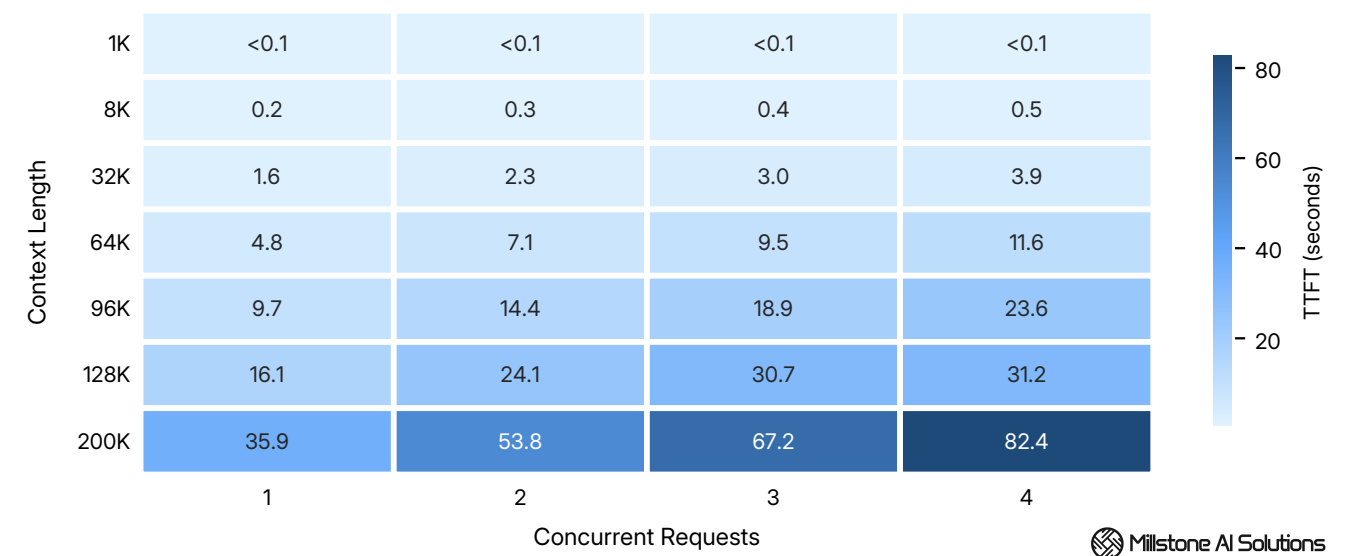
LATENCY

Time to First Token

Time from request submission to first response token. The primary metric for perceived responsiveness. TTFT has two components:

- **Queue wait:** Time waiting for GPU availability (increases with concurrency)
- **Prefill:** Time to process input context (increases with context length)

At low concurrency, prefill dominates. Under load, queue wait takes over. See Technical Analysis for more.



How Responsive is This?

TTFT	USER EXPERIENCE
< 500ms	Feels instant
500ms–2s	Feels responsive
2–5s	Noticeable but still acceptable
5–10s	Feels slow; generally acceptable at higher context lengths
> 10s	Can be frustrating; users may retry or abandon

Important note about caching. These benchmarks use fresh context with no caching enabled, representing worst-case TTFT. In production with caching enabled, only new tokens require processing. For example, a 64K conversation where you add 1K of new context would have a TTFT similar to the 1K results above, not the 64K results. **For most real-world use cases where context is built incrementally (chatbots, coding assistants, multi-turn agents), TTFT with caching enabled would be significantly faster than these results.**

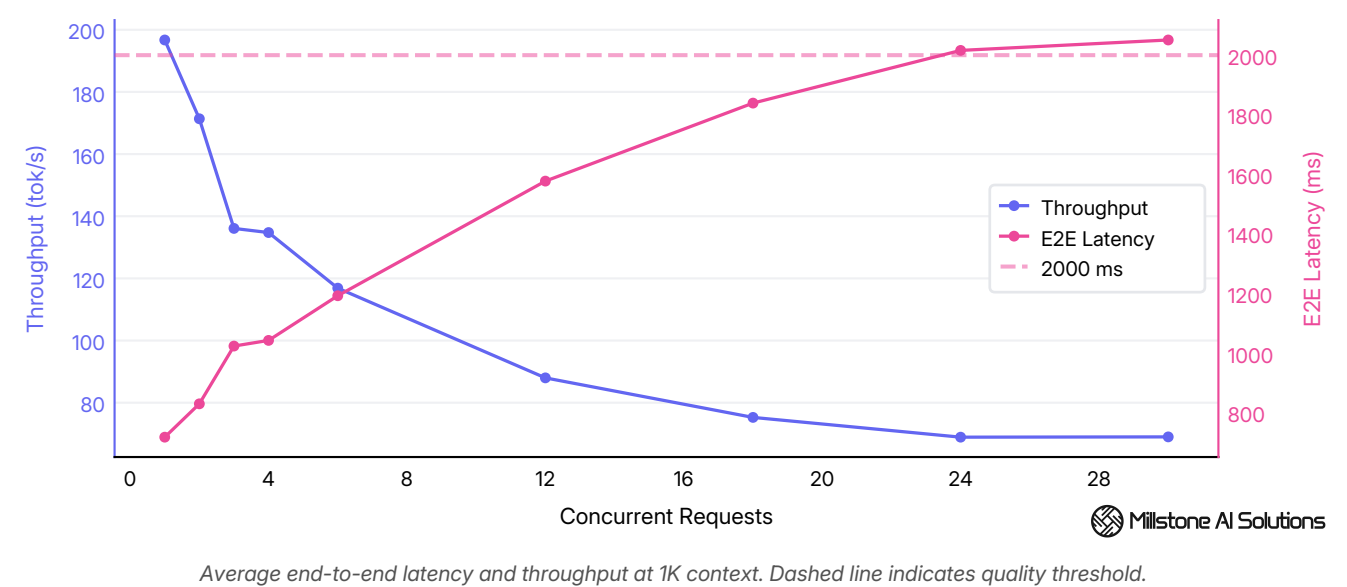
Capacity Analysis

How many concurrent requests can this configuration handle for different workloads? Each section below shows performance metrics as concurrency increases at a specific context length. Dashed lines indicate quality thresholds, the point where user experience degrades below acceptable levels. The "capacity limit" is the tested or estimated point where the first threshold is reached.

Code Completion (1K Context)

Inline code suggestions in IDEs, like autocomplete. Responsiveness is critical. This test generates 128 output tokens per request (vs. 1024 elsewhere) to match typical autocomplete length. The key metric is end-to-end latency, not TTFT.

Threshold: End-to-end latency < 2,000ms



METRIC	@ 1 request	@ 24 requests	@ 30 requests
End-to-end latency	717ms	2016ms (threshold exceeded)	2051ms (threshold exceeded)
Throughput	197 tok/s	69 tok/s	69 tok/s

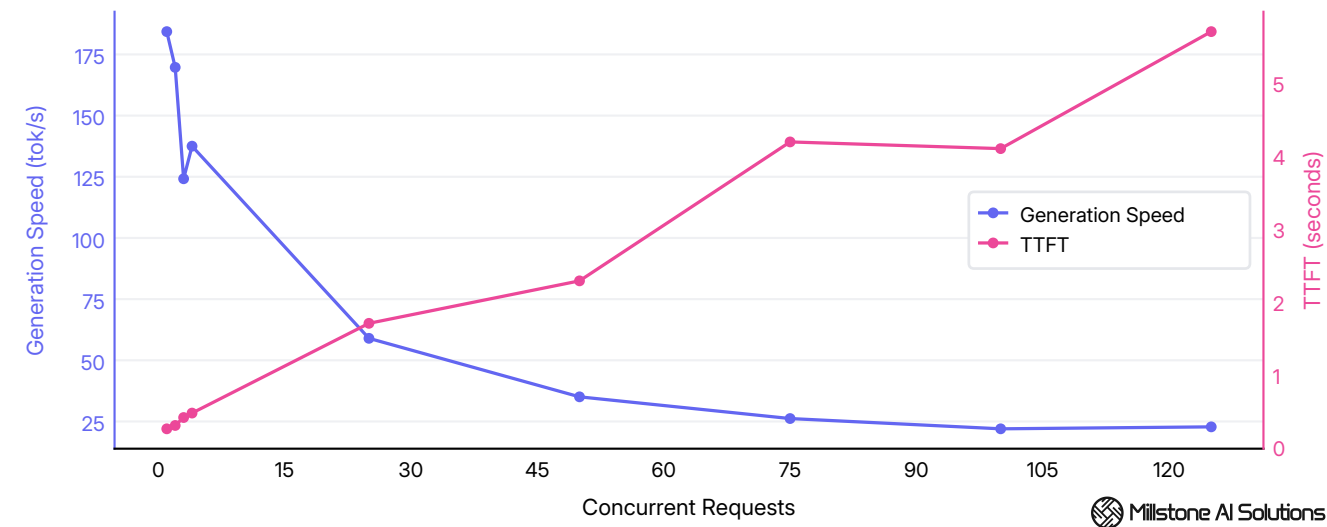
Capacity limit: 24 concurrent requests

At 24 concurrent requests, end-to-end latency reaches 2016ms, just above the 2,000ms threshold.

Short-form Chatbot (8K Context)

Quick conversational exchanges: customer support queries, simple Q&A, single-turn requests. 8K context accommodates a few back-and-forth messages plus system prompt. User expectations are more forgiving for these scenarios. 10+ tok/s is acceptable for reading streamed responses from a support chatbot.

Thresholds: TTFT < 10s, generation speed > 10 tok/s



Average per-user generation speed and TTFT at 8K context.

METRIC	@ 1 request	@ 75 requests	@ 125 requests
TTFT	0.2s	4.2s	5.7s
Generation speed	184 tok/s	26 tok/s	23 tok/s

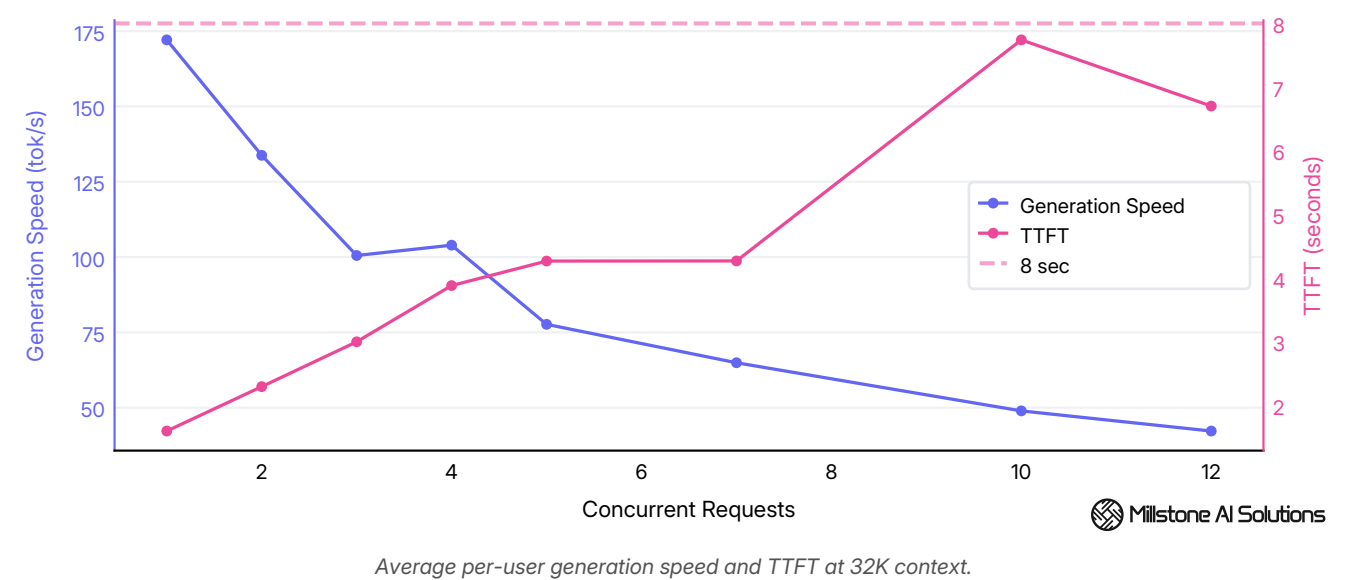
Capacity limit: 125+ concurrent requests

At 125 concurrent requests, TTFT is 5.7 seconds and generation speed is 23 tok/s, both well within acceptable bounds. Capacity likely extends higher.

General Chatbot (32K Context)

ChatGPT-style chatbot. If you're deploying a multi-turn conversational chatbot, this benchmark shows how many concurrent requests you can support while matching acceptable responsiveness. 32K context matches ChatGPT's limit.

Thresholds: TTFT < 8s, generation speed > 15 tok/s



METRIC	@ 1 request	@ 5 requests	@ 12 requests
TTFT	1.6s	4.3s	6.7s
Generation speed	172 tok/s	78 tok/s	42 tok/s

Capacity limit: 12+ concurrent requests

At 12 concurrent requests, TTFT is 6.7 seconds and generation speed is 42 tok/s, both well within acceptable bounds. The configuration handles this workload comfortably within tested limits; capacity likely extends higher.

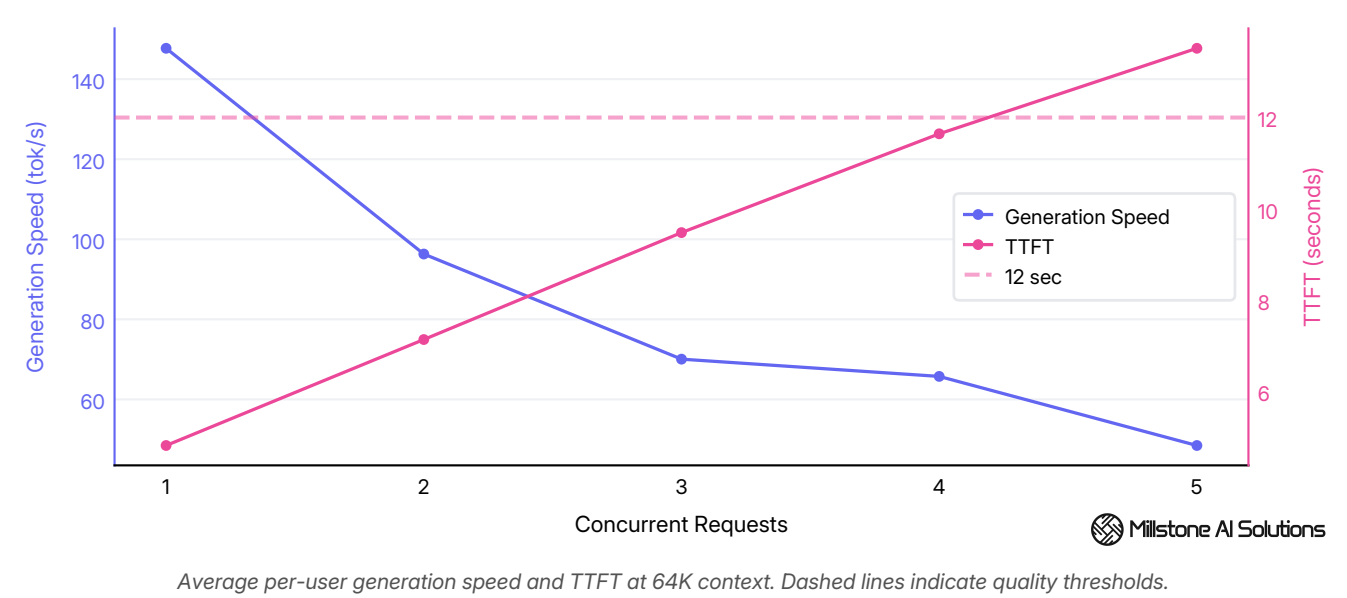
Note about caching: Most chatbot users build context incrementally over a conversation. With caching properly configured, TTFT is dramatically reduced since only new tokens require processing. These results represent worst-case TTFT where all context is processed at once.

Long Document Processing (64K Context)

Summarizing reports, extracting data from contracts, analyzing lengthy documents. 64K tokens handles documents up to roughly 125-160 pages depending on formatting and density.

Users typically tolerate higher latency for document processing since they understand large inputs require more processing time. However, generation speed still needs to stay at or above reading speed.

Thresholds: TTFT < 12s, generation speed > 15 tok/s



METRIC	@ 1 request	@ 4 requests	@ 5 requests
TTFT	4.8s	11.6s	13.5s (threshold exceeded)
Generation speed	148 tok/s	66 tok/s	48 tok/s

Capacity limit: 4 concurrent requests

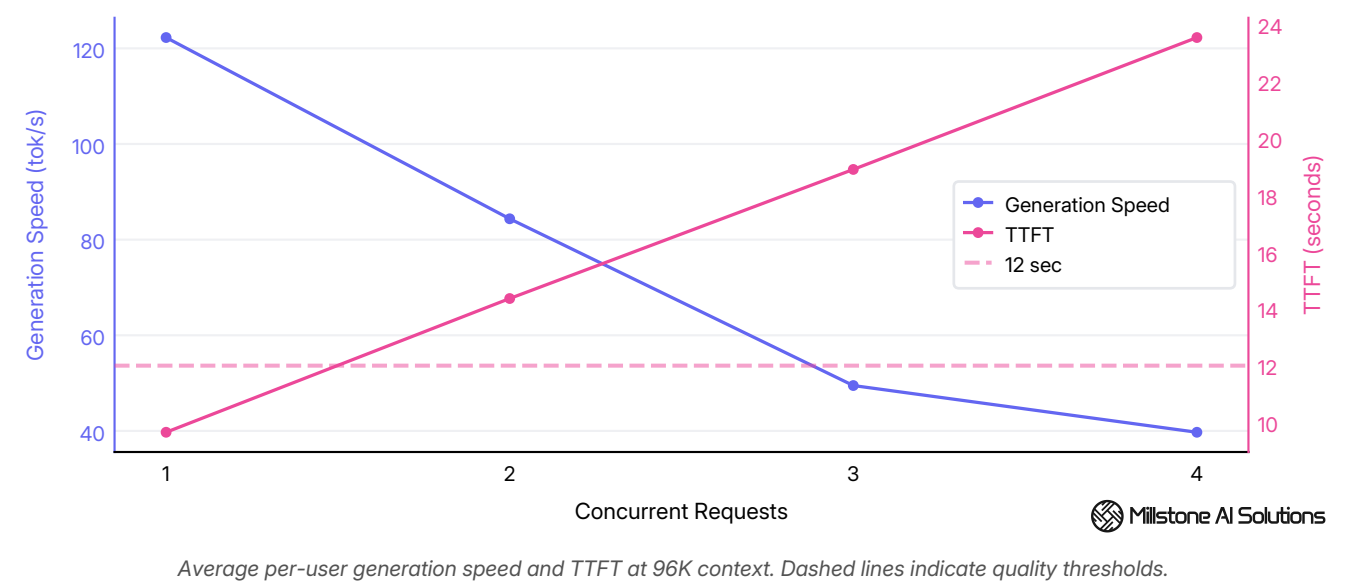
At 4 concurrent requests, TTFT reaches 11.6 seconds, just under the 12-second threshold. Generation speed at this concurrency is 66 tok/s, above the 15 tok/s minimum.

Automated Coding Assistant (96K Context)

Agentic coding workloads: AI assistants that read large portions of a codebase to answer questions, refactor code, or implement features. 96K tokens handles roughly 8,000-9,000 lines of code, enough for significant repository context.

Agentic workflows chain multiple LLM calls (tool use, retrieval, iterative refinement). With caching properly configured, context persists between requests and only new tokens require processing, dramatically reducing TTFT for each step. These results represent worst-case TTFT where all context is processed at once.

Thresholds: TTFT < 12s, generation speed > 20 tok/s



METRIC	@ 1 request	@ 2 requests	@ 4 requests
TTFT	9.7s	14.4s (threshold exceeded)	23.6s (threshold exceeded)
Generation speed	122 tok/s	84 tok/s	40 tok/s

Capacity limit: 1 request

This configuration handles single-user agentic coding workloads with 9.7s TTFT and 122 tok/s generation speed, acceptable for individual use.

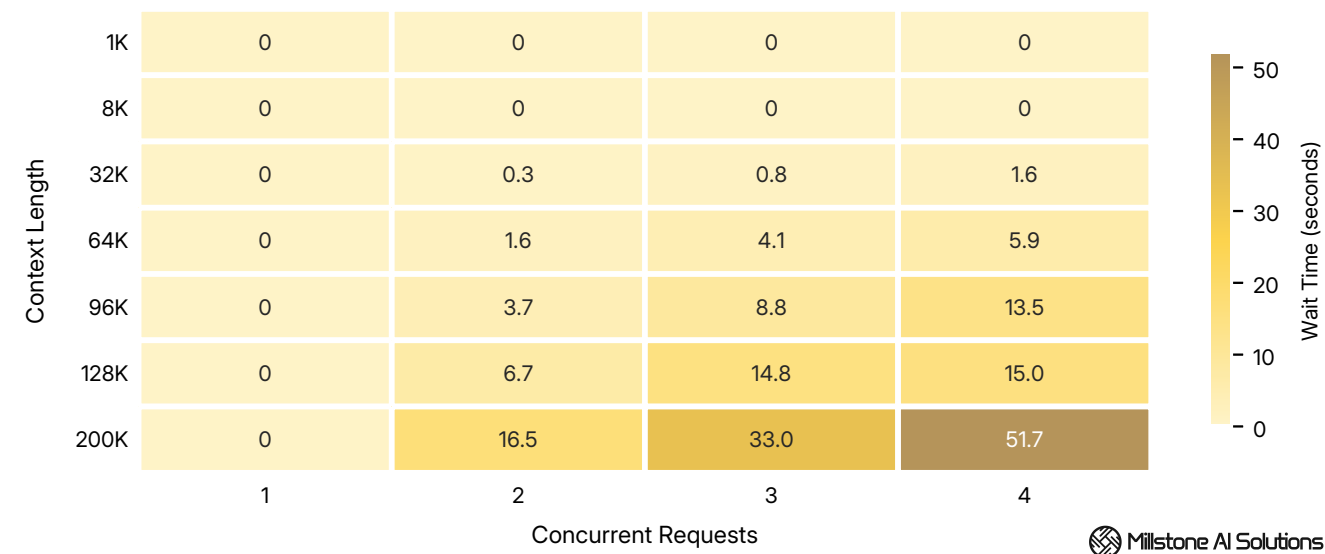
Technical Analysis

Infrastructure-level metrics that explain user-facing performance. Queue depth, prefill throughput, token generation latency, and scaling efficiency across load conditions. These help diagnose bottlenecks and validate infrastructure decisions.

Queue Wait Times

Time a request waits for GPU availability before processing begins. At low concurrency, queue wait is near zero. As load increases, requests queue and wait times grow.

Queue wait is included in TTFT. Breaking it out separately helps identify whether latency is caused by GPU saturation (high queue wait) or context processing (high prefill time).



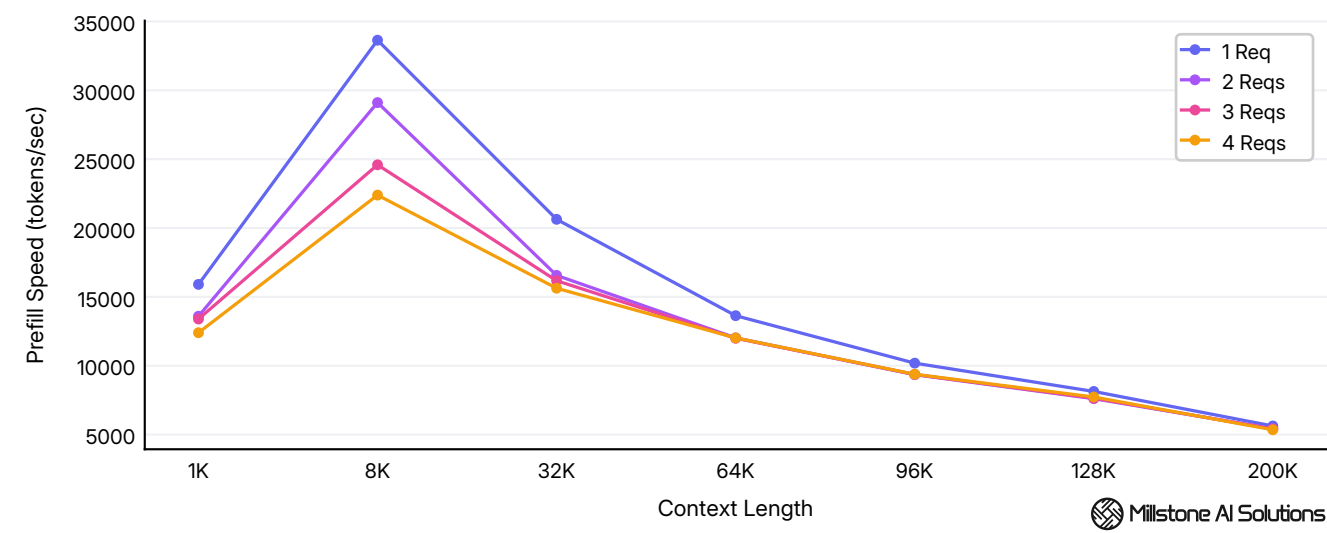
Average queue wait time across 1K - 200K tokens context at 1 - 4 concurrent requests.

At single concurrency, queue wait is effectively zero regardless of context length. At **4 concurrent requests** with **200K context**, queue wait reaches **51.7 seconds**. Rising queue times signal GPU saturation, meaning requests are waiting for resources rather than being processed immediately.

Interpretation: Queue wait time and prefill time are measured independently and may not sum exactly to TTFT. Under heavy load, chunked prefill and preemptions can cause these metrics to overlap, sometimes resulting in queue wait + prefill exceeding TTFT. Use queue wait for capacity planning and identifying bottlenecks. Use TTFT for actual user wait time before streaming begins.

Per-User Prefill Speed

Rate at which the model processes input context before generating output. Prefill speed determines the non-queue portion of TTFT. Higher prefill speeds mean faster time-to-first-token at a given context length.



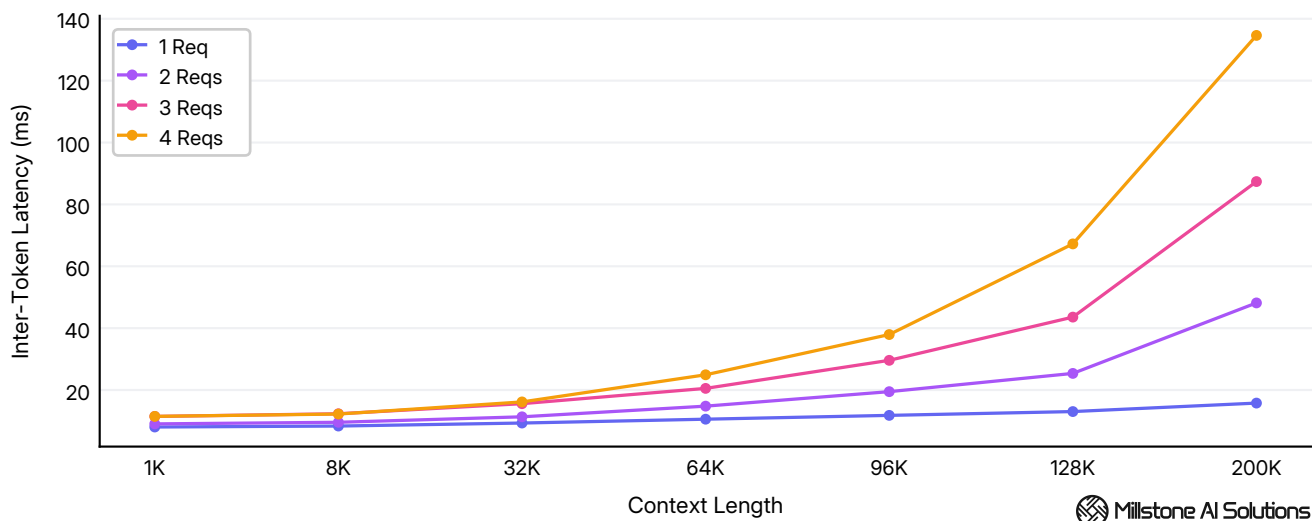
Average per-user prefill speed across 1K - 200K tokens context at 1 - 4 concurrent requests.

CONCURRENT REQUESTS	PEAKS AT	PEAK SPEED
1	8K context	33,639 tok/s
2	8K context	29,109 tok/s
3	8K context	24,591 tok/s
4	8K context	22,391 tok/s

Prefill speed peaks at a certain context length and then declines as additional context increases computational overhead. This peak can reflect GPU saturation (compute or memory bandwidth fully utilized) or engine configuration such as chunked prefill limits, which cap tokens processed per forward pass to maintain responsiveness under load. On the chart, this appears as lines that peak before reaching the longest context.

Inter-Token Latency

Time between consecutive tokens during generation. Determines the smoothness of responses. Lower latency produces more fluid output. ITL helps diagnose the underlying token-level behavior.

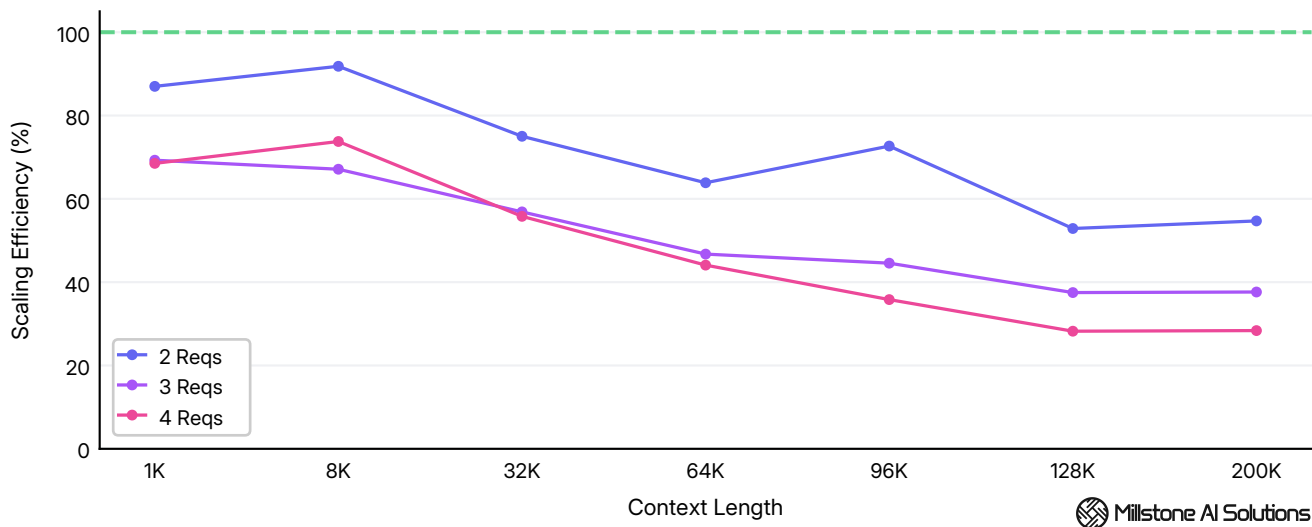


Average inter-token latency across 1K - 200K tokens context at 1 - 4 concurrent requests.

At single-user short context, inter-token latency is imperceptible (**8ms**). The highest latency observed was **135ms** at **200K** context with **4 concurrent requests**, where individual tokens become visible as they stream.

Scaling Efficiency

Percentage of ideal linear scaling achieved as concurrency increases. 100% efficiency means doubling concurrent requests doubles total throughput with no per-user degradation. Real-world efficiency is always lower due to shared GPU resources.



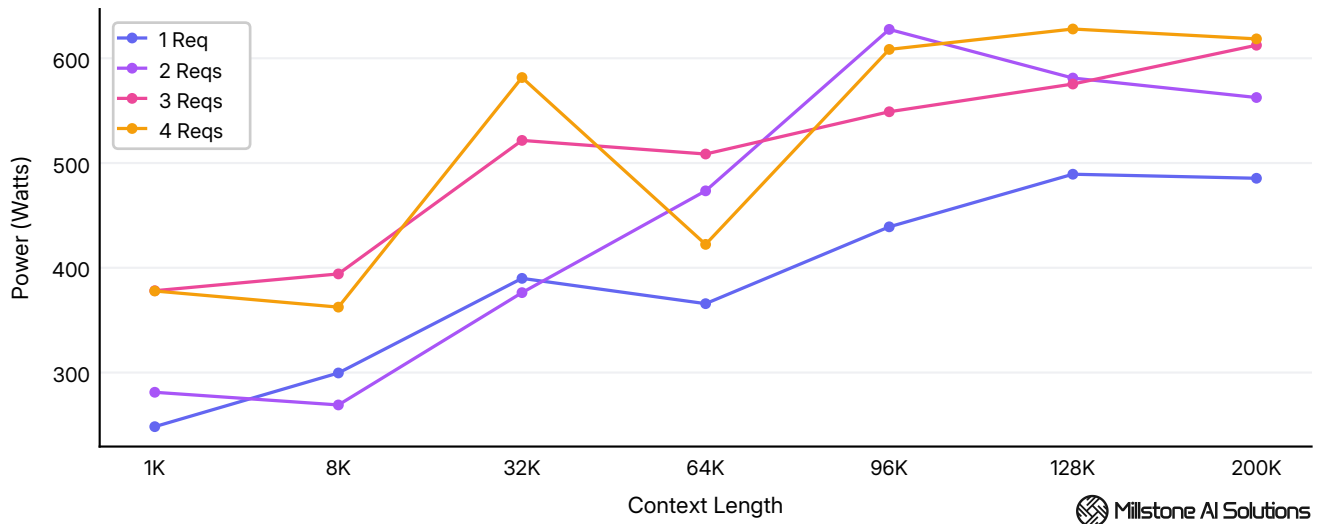
Scaling efficiency across 1K - 200K tokens context at 1 - 4 concurrent requests.

Efficiency remains high at low concurrency where GPU resources can serve requests without contention. At higher concurrency, efficiency drops as requests compete for shared resources. High efficiency at your target concurrency indicates headroom for growth. Sharply dropping efficiency signals diminishing returns.

EFFICIENCY

Power Consumption

GPU power draw under varying load conditions. Relevant for operational cost estimation, cooling requirements, and data center power budgeting.



Average GPU power draw across 1K - 200K tokens context at 1 - 4 concurrent requests.

Power consumption scales with both context length and concurrency. The highest power draw observed was **628W** at **128K** context with **4 concurrent requests**, costing approximately **\$0.06/hour** at \$0.10/kWh. Higher concurrency or sustained load beyond tested conditions may increase power consumption further. For infrastructure planning, budget for peak power draw.

Need Help Deciding?

Not sure what configuration you need? Our team can help you identify the right model, hardware, and deployment strategy for your specific use case.

[Schedule a Conversation →](#)

Additional data available on request: full percentile breakdowns (P50–P99) and GPU metrics.