

Technical Report

# Moreh vLLM Performance Evaluation: DeepSeek V3/R1 671B on AMD Instinct MI300X GPUs

Moreh, Inc.

August 2025

# Contents

Overview .....	1
AMD Instinct MI300X GPU .....	2
Optimizations for DeepSeek V3/R1 671B .....	3
Experimental Setup.....	3
Output TPS, TTFT, and TPOT .....	4
Trade-Off Between Latency and Throughput.....	7
Conclusion .....	7
Appendix: Raw Data.....	8

# Moreh vLLM Performance Evaluation: DeepSeek V3/R1 671B on AMD Instinct MI300X GPUs

## Overview

Moreh develops software to enable various AI workloads – from pretraining to inference – to run efficiently on non-NVIDIA accelerators, with a particular focus on AMD GPUs.

vLLM is one of the most widely adopted inference engines for running LLM services in research, enterprise, and production environments. It is developed by a strong open-source community with contributions from both academia and industry, and provides broad support for various models, hardware, and optimization techniques. AMD is also contributing to the project to make vLLM run on AMD GPUs and the ROCm software stack. Nevertheless, most optimizations in vLLM still target NVIDIA GPUs, and the performance of AMD GPU hardware has yet to be fully utilized.

Moreh vLLM is our optimized version of vLLM, designed to deliver superior LLM inference performance on AMD GPUs. It supports the same models and features as the original vLLM, while maximizing computational performance on the AMD CDNA architecture. This is achieved through Moreh's proprietary compute and communication libraries, along with model-level optimizations and vLLM engine-level modifications.

This technical report evaluates the inference performance of the DeepSeek V3/R1 671B model – one of the most advanced open-source LLMs available today – on Moreh vLLM. We conduct comprehensive testing across various input/output lengths and concurrency levels. Compared to the original vLLM, Moreh vLLM delivers an average of **1.68x** higher throughput (total output tokens per second). Furthermore, it reduces latency metrics (time to first token and time per output token) by an average of **1.75x** and **1.70x**, respectively. In conclusion, adopting Moreh vLLM unlocks the full potential of AMD MI300 series GPUs, enabling them to serve as an efficient inference system.

## AMD Instinct MI300X GPU

The AMD Instinct MI300X GPU presents a compelling alternative to NVIDIA's H100. It provides 1.32x higher theoretical compute performance, 2.4x larger memory capacity, and 1.58x higher peak memory bandwidth compared to the H100. In particular, its significantly larger memory capacity and bandwidth are a major advantage for optimizing LLM inference. Table 1 compares the detailed hardware specifications.

Table 1. Comparison between NVIDIA H100 and AMD MI300X

Items	H100 SXM	MI300X	Relative (MI300X/H100)
<b>Basic facts</b>			
Architecture	Hopper	CDNA3	
Form factor	SXM5 module	OAM module	
Lithography	TSMC 4 nm	TSMC 5 nm	
# SMs (compute units)	132	304	
# cores	16,896	19,456	
# tensor/matrix cores	528	1,216	
Peak engine clock	1,830 MHz	2,100 MHz	
TDP	700 W	750 W	
<b>Peak theoretical performance (dense)</b>			
FP32 vector	66.9 TFLOPS	163.4 TFLOPS	2.44x
TF32 matrix	494.7 TFLOPS	653.7 TFLOPS	1.32x
FP16/BF16 matrix	989.4 TFLOPS	1,307.4 TFLOPS	1.32x
FP8 matrix	1978.9 TFLOPS	2,614.9 TFLOPS	1.32x
INT8 matrix	1978.9 TOPS	2,614.9 TOPS	1.32x
<b>GPU memory</b>			
Technology	HBM3	HBM3	-
Capacity	80 GB	192 GB	2.40x
Peak bandwidth	3.35 TB/s	5.3 TB/s	1.58x
<b>Cache and scratchpad</b>			
L1D + scratchpad	256 KB per SM	32+64 KB per SM	0.38x
L2/L3	50 MB L2	32 MB L2, 256 MB L3	-
<b>Connectivity (H2D: host to device, D2D: device to device within a server)</b>			
H2D interface	PCIe Gen5 x16	PCIe Gen5 x16	-
H2D bandwidth	128 GB/s	128 GB/s	1.00x
D2D interface	NVLink Gen4	Infinity Fabric Gen4	-
D2D bandwidth	900 GB/s	896 GB/s	0.996x
# GPUs per server	8	8	1.00x

AMD has also released the MI325X and MI355X as successors to the MI300X, which are direct competitors to NVIDIA's H200 and B200 GPUs, respectively. Since these next-generation models are also based on the AMD CDNA3 architecture, all optimizations within Moreh vLLM will continue to apply seamlessly. We plan to publish performance evaluation results on the MI325X and MI355X in the near future and are always open to partners who can provide development and testing servers.

## Optimizations for DeepSeek V3/R1 671B

Moreh vLLM incorporates numerous optimizations to enhance the performance of the DeepSeek 671B model, including, but not limited to:

- **Optimal GEMM and Attention Kernel Selection:** To achieve consistently high performance across various scenarios (e.g., different input/output sequence lengths and batch sizes), Moreh vLLM dynamically selects the optimal GEMM and Attention kernels without the need for online profiling and manual tuning.
- **Fused MoE Kernel Optimization:** We have implemented a highly optimized fused MoE kernel that delivers better performance than AMD's AITER library, particularly for small batch sizes.
- **FP8 KV Cache Support:** Moreh vLLM includes Mult-head Latent Attention (MLA) kernels that enables the KV cache to be stored and loaded in FP8 format. This optimization significantly improves performance, especially in long-context scenarios.
- **Vertical and Horizontal Kernel Fusion:** Moreh vLLM employs both vertical fusion (e.g., fused RoPE kernels) and horizontal fusion (e.g., merging multiple GEMMs in shared experts) to reduce kernel launch overhead and improve computational efficiency.
- **vLLM Engine-Level Modifications:** We have made modifications at the vLLM engine level to more efficiently utilize AMD GPUs, including leveraging HIP graphs for streamlined kernel execution.

## Experimental Setup

All experiments were conducted on an MI300X server configured as follows:

- **Server:** Lenovo ThinkSystem SR685a V3
- **CPU:** 2x AMD EPYC 9534 (128 cores in total, 2.45 GHz)
- **GPU:** 8x AMD Instinct MI300X OAM
- **Main Memory:** 2,304 GB (24x 96 GB)
- **Operating System:** Ubuntu 22.04.4 (Linux kernel 5.15.0-25-generic)
- **ROCm Version:** 6.8.5

We used the open-source vLLM 0.9.2 (tag v0.9.2 of <https://github.com/ROCm/vllm>) as a baseline for comparison. This was the latest versions available at the time of testing.

The DeepSeek model was executed in parallel across 8 GPUs of the server with a tensor parallelism (TP) of 8. Thanks to AMD MI300X's large memory capacity of 192 GB, over half of the GPU memory remains available even after storing ~84 billion parameters per GPU in FP8 format. This allows the server to handle numerous requests with high concurrency, showcasing a significant advantage for large-scale generative AI workloads.

Performance was measured using vLLM's benchmark\_serving tool. We chose 70 different combinations of input sequence length (ISL), output sequence length (OSL), and concurrency, as shown in Table 2.

The experimental setup was determined through discussions with one of our customers in Korea.

Table 2. Various request patterns used for performance measurement

Input sequence length (ISL)	Output sequence length (OSL)	Concurrencies
1024	1024	1, 2, 4, 8, 16, 32, 64, 128, 256, 512
1024	4096	1, 2, 4, 8, 16, 32, 64, 128, 256, 512
4096	1024	1, 2, 4, 8, 16, 32, 64, 128, 256, 512
4096	4096	1, 2, 4, 8, 16, 32, 64, 128, 256, 512
16384	1024	1, 2, 4, 8, 16, 32, 64, 128
16384	4096	1, 2, 4, 8, 16, 32, 64, 128
32768	1024	1, 2, 4, 8, 16, 32, 64
32768	4096	1, 2, 4, 8, 16, 32, 64

## Output TPS, TTFT, and TPOT

Output tokens per second (TPS), time to first token (TTFT), and time per output token (TPOT) are three key metrics for evaluating the performance of LLM inference.

- *Output tokens per second* measures the overall throughput of the system, indicating how many tokens the model can generate in one second across all concurrent requests.
- *Time to first token* captures the initial latency – the time from when a request is sent until the very first token is produced.
- *Time per output token* indicates the average time taken to generate each subsequent token after the first one.

Output tokens per second is directly tied to service cost (dollar per token). The latter two metrics are important for user-perceived responsiveness. Together, measuring these three metrics provides a comprehensive view of inference performance, balancing cost and user experience.

Figure 1 shows a graph comparing output tokens per second. Figure 2 and Figure 3 present graphs comparing the mean time to first token and the mean time per output token, respectively. The raw data can be found in the appendix. Moreh vLLM achieves 1.68x higher total output tokens per second, 1.75x lower time to first token, and 1.7x lower time per output token compared to the original vLLM. This demonstrates that simply replacing the software with Moreh vLLM on the same AMD MI300 series GPU system can reduce costs while improving user experience.

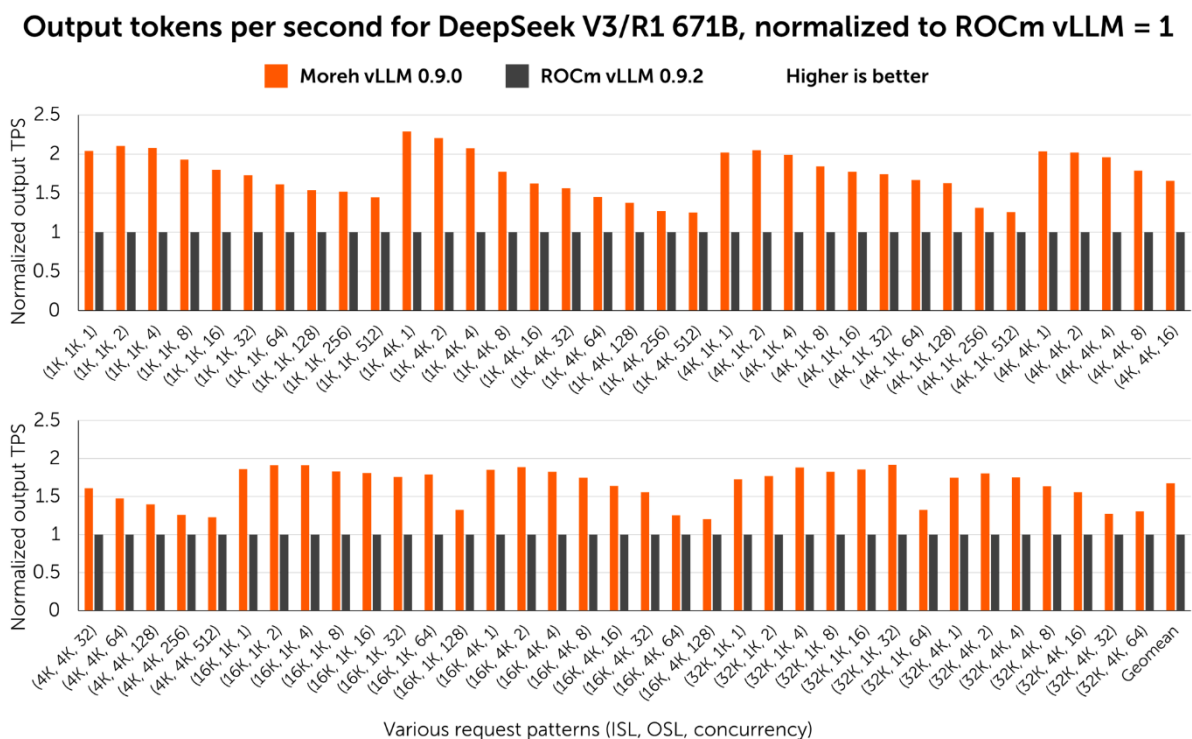


Figure 1. Output tokens per second for various request patterns. Higher is better. Moreh vLLM shows an average of 1.68x higher performance.

### Mean time to first token for DeepSeek V3/R1 671B, normalized to ROCm vLLM = 1

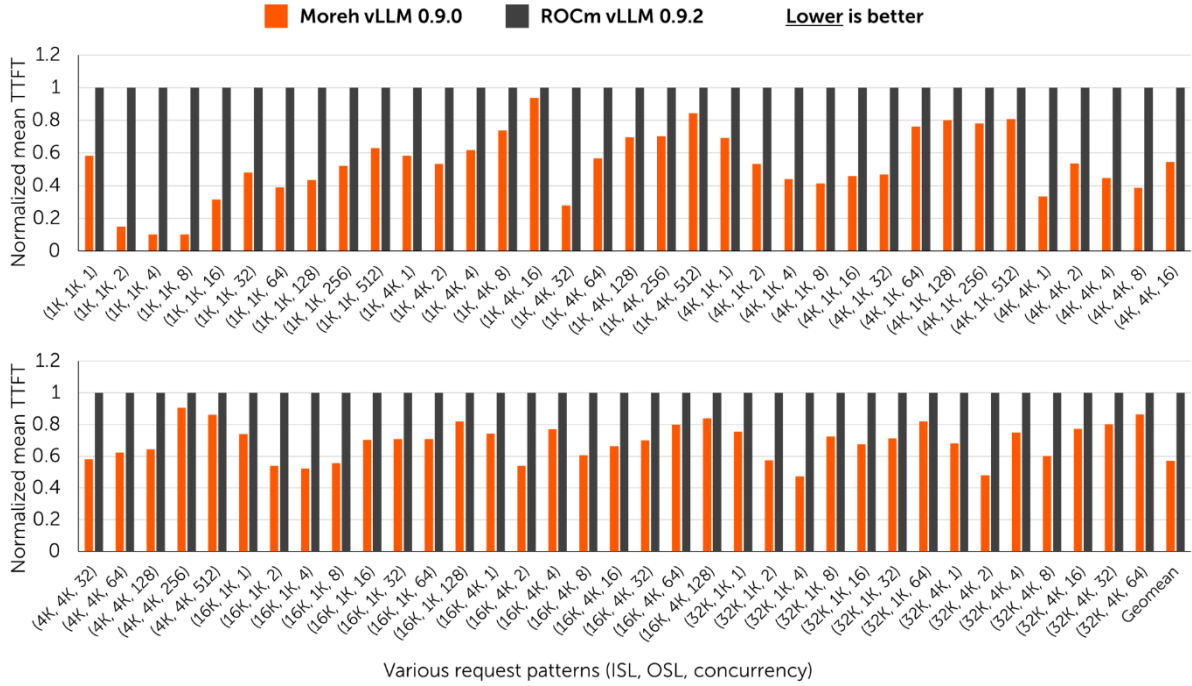


Figure 2. Mean time to first token for various request patterns. Lower is better. Moreh vLLM shows an average of 1.75x lower latency.

### Mean time per output token for DeepSeek V3/R1 671B, normalized to ROCm vLLM = 1

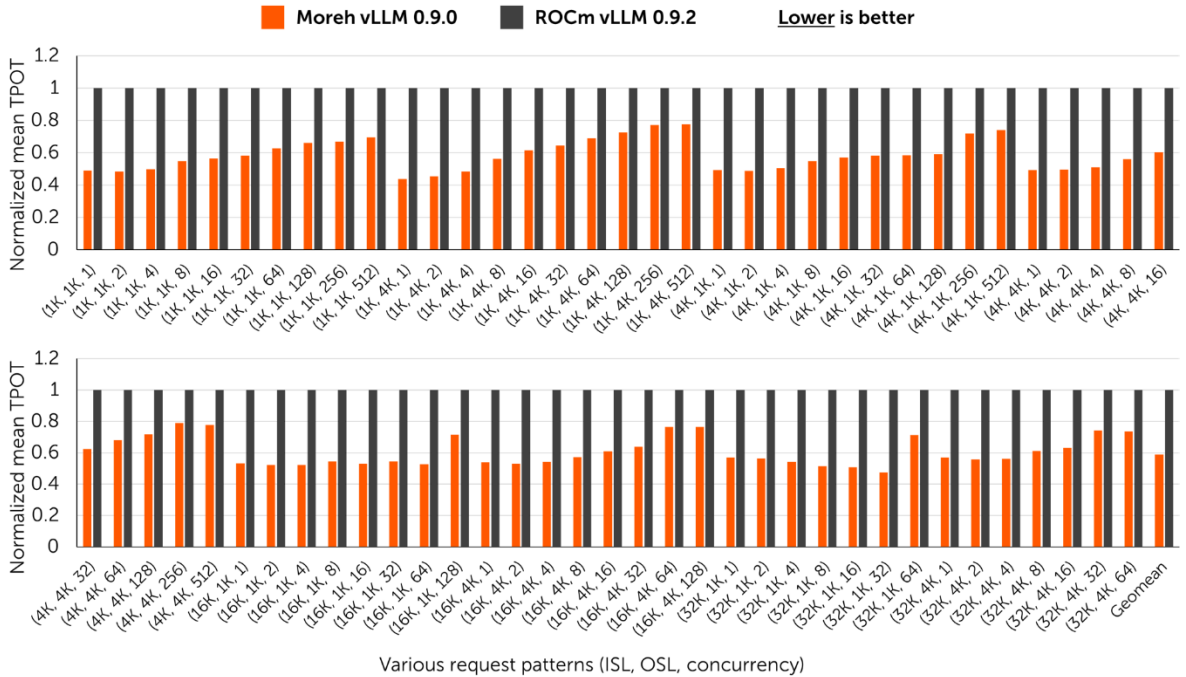


Figure 3. Mean time per output token for various request patterns. Lower is better. Moreh vLLM shows an average of 1.70x lower latency.



## Trade-Off Between Latency and Throughput

LLM inference involves an inherent trade-off between latency and throughput. Increasing the maximum concurrency of a vLLM instance improves throughput but also increases latency, while decreasing concurrency improves latency but lowers throughput.

Figure 4 illustrates these latency-throughput trade-off curves for the original vLLM and Moreh vLLM across various request patterns (input/output sequence lengths). Overall, the closer the graph shifts toward the upper left, the better the performance characteristics.

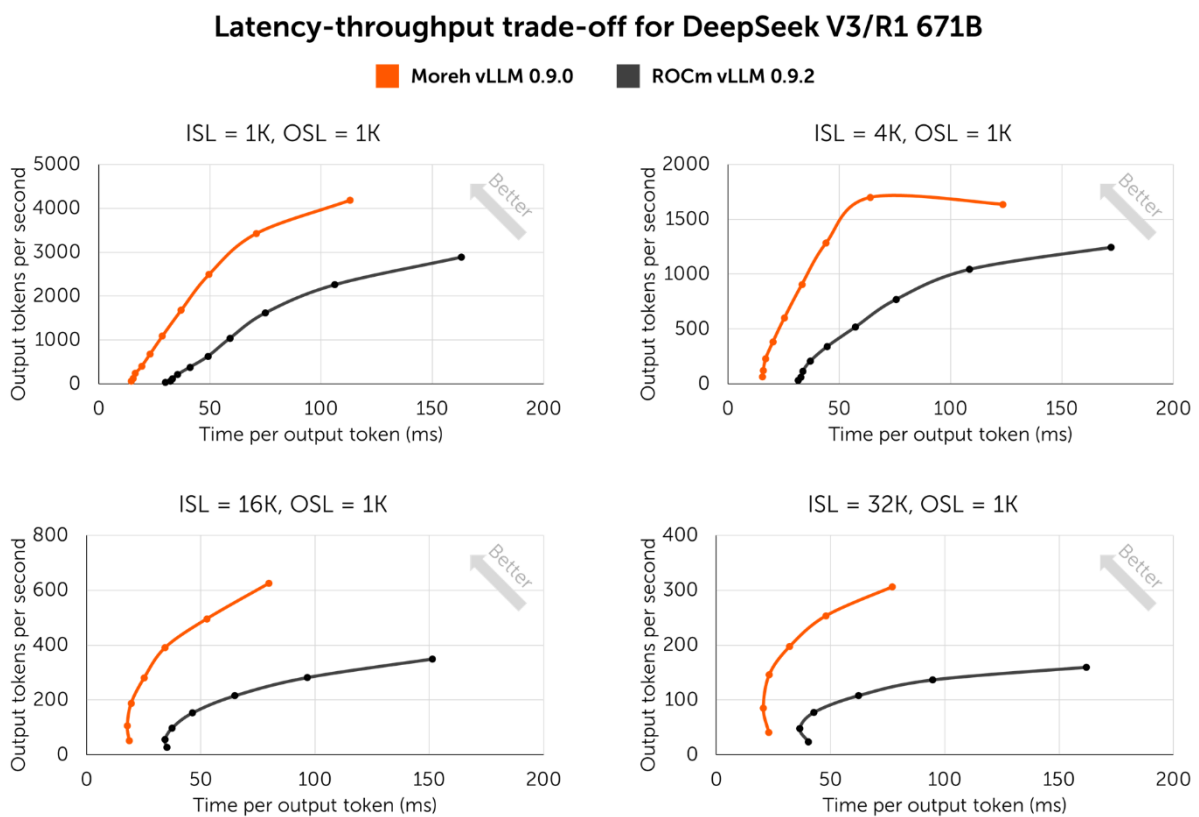


Figure 4. Trade-off curves between time per output token (latency) and output tokens per second (throughput), for different input/output sequence lengths.

## Conclusion

Moreh vLLM incorporates various techniques to optimize inference for the DeepSeek V3/R1 model, including proprietary GPU libraries, model-level optimizations, and modifications to the vLLM engine. As a result, Moreh vLLM delivers substantial performance improvements over the original open-source vLLM across various inference metrics. By adopting Moreh vLLM on AMD MI300 series GPU servers, LLM

services can reduce costs while simultaneously improving latency. Moreh also provides a service that optimizes a customer’s proprietary AI model on AMD GPUs and delivers on-demand vLLM for it.

## Appendix: Raw Data

Request patterns			Moreh vLLM 0.9.0			ROCm vLLM 0.9.2		
ISL	OSL	Conc.	Output TPS	Mean TTFT (ms)	Mean TPOT (ms)	Output TPS	Mean TTFT (ms)	Mean TPOT (ms)
1024	1024	1	67.66	68.87	14.73	33.19	117.75	30.05
1024	1024	2	127.42	101.03	15.61	60.66	673.78	32.34
1024	1024	4	240.54	157.81	16.49	115.68	1570.5	33.07
1024	1024	8	407.51	258.76	19.39	211.3	2576.59	35.37
1024	1024	16	681.79	371.38	23.12	379.63	1178.78	41.01
1024	1024	32	1088.31	789.45	28.63	630	1640.65	49.19
1024	1024	64	1681.08	946.86	37.12	1041.09	2434.3	59.06
1024	1024	128	2491.89	1738.42	49.56	1618.89	3985.35	74.99
1024	1024	256	3426.59	3625.25	70.81	2258.75	6942.03	106.09
1024	1024	512	4188.01	8328.91	113.07	2888.9	13220.19	163.02
1024	4096	1	67.84	26.25	14.74	29.66	44.9	33.72
1024	4096	2	128.8	36.97	15.52	58.48	69.19	34.19
1024	4096	4	241.82	49.13	16.53	116.79	79.56	34.24
1024	4096	8	408.53	68.11	19.57	230.29	92.15	34.72
1024	4096	16	668.38	115.34	23.91	412.03	123.06	38.81
1024	4096	32	1050.95	695.65	30.28	671.76	2496.33	47.02
1024	4096	64	1647.27	1219.32	38.54	1132.97	2151.3	55.95
1024	4096	128	2400.07	2469.2	52.68	1741.15	3544.58	72.6
1024	4096	256	2989.84	8761.55	79.91	2347.09	12489.49	103.73
1024	4096	512	2865.59	50932.74	140.1	2291.7	60350.75	180.37
4096	1024	1	63.56	168.89	15.58	31.5	244.11	31.54
4096	1024	2	123.53	241.76	15.96	60.31	452.43	32.74
4096	1024	4	230.56	382.56	16.98	116.04	871.03	33.62
4096	1024	8	383.13	641.39	20.25	207.82	1555.13	36.97
4096	1024	16	602.42	1148.78	25.41	340.33	2501.4	44.54
4096	1024	32	905.92	2037.01	33.24	519.28	4351.16	57.25
4096	1024	64	1283.34	5667.4	44.09	770.49	7444.9	75.48
4096	1024	128	1699.79	10904.04	64.01	1044.22	13633.68	108.44
4096	1024	256	1636.3	21452.36	123.44	1245.4	27480.28	172.01
4096	1024	512	1498.15	145456.58	152.69	1191.79	180087.52	206.4
4096	4096	1	62.02	169.1	16.09	30.49	504.59	32.69
4096	4096	2	122.22	246.52	16.31	60.53	460.12	32.94
4096	4096	4	231.21	389.42	17.21	117.97	873.01	33.7
4096	4096	8	389.51	643.49	20.38	217.6	1666.97	36.36
4096	4096	16	619.44	1344.41	25.49	373.34	2462.39	42.24
4096	4096	32	972.65	2406.18	32.28	605.53	4135.62	51.8
4096	4096	64	1426.07	4568.72	43.69	968.33	7340.71	64.21
4096	4096	128	1955.28	8842.49	63.1	1398.04	13732.72	87.97
4096	4096	256	1736.1	85503.85	109.68	1380.32	94500.17	138.72

4096	4096	512	1654.59	515697.59	128.44	1351.04	599442.6	165.24
16384	1024	1	51.55	679.9	18.75	27.71	918.52	35.23
16384	1024	2	106.01	988.41	17.9	55.48	1828.67	34.28
16384	1024	4	187.84	1749.99	19.56	98.23	3354.89	37.42
16384	1024	8	280.33	3313	25.23	153.12	5958.53	46.33
16384	1024	16	390.64	6513.2	34.41	215.59	9260.04	64.92
16384	1024	32	495.8	11657.99	52.68	282.09	16496.45	96.69
16384	1024	64	624.79	21913.35	79.83	349.39	30946.86	151.37
16384	1024	128	435.75	130014.42	138.95	328.6	158813.8	194.41
16384	4096	1	51.87	680.16	19.12	28.04	916.8	35.45
16384	4096	2	110.32	985.81	17.89	58.4	1825.71	33.8
16384	4096	4	202.79	2474.91	19.11	111.04	3218.31	35.23
16384	4096	8	328	3308.94	23.56	187.79	5459.98	41.24
16384	4096	16	495.07	6029.75	30.79	302.31	9106.35	50.63
16384	4096	32	676.93	11445.82	44.35	435.15	16348.36	69.37
16384	4096	64	682.57	26408.54	78.04	546.01	33073.29	102.09
16384	4096	128	633.63	312895.79	94.46	526.67	372806.36	123.42
32768	1024	1	40.76	1619.06	22.98	23.59	2142.13	40.35
32768	1024	2	85.07	2947.79	20.62	48.07	5131.4	36.58
32768	1024	4	145.97	4265.61	23.17	77.46	9006.42	42.76
32768	1024	8	197.27	8546.05	32.03	108.06	11793.94	62.28
32768	1024	16	253.62	15032.53	47.97	136.75	22232.48	94.7
32768	1024	32	306.4	26950.98	77.01	159.79	37871.35	161.83
32768	1024	64	197.15	145957.44	154.46	148.77	178245.2	216.6
32768	4096	1	42.23	1629.71	23.29	24.14	2398.39	40.84
32768	4096	2	97.48	2449.27	19.91	53.99	5101.95	35.79
32768	4096	4	173.55	5263.17	21.74	98.91	7008.72	38.71
32768	4096	8	264.33	8294.78	28.19	161.55	13814.27	46.08
32768	4096	16	363.12	15964.79	40.03	233.25	20647.09	63.4
32768	4096	32	396.86	29920.93	69.14	312.04	37258.33	93.07
32768	4096	64	338.29	318451.11	93.56	259.36	368561.16	127.47



To learn more, please visit our website (<https://moreh.io>) or contact us ([contact@moreh.io](mailto:contact@moreh.io)).

Copyright ©2025 Moreh, Inc. All rights reserved.

The information contained herein is for informational purposes only and is subject to change without notice. While every precaution has been taken in the preparation of this document, it may contain technical inaccuracies, omissions, and typographical errors, and Moreh, Inc. is under no obligation to update or otherwise correct this information. Moreh, Inc. makes no representations or warranties with respect to the accuracy or completeness of the contents of this document and assumes no liability of any kind for the consequences or use of such information or for any infringement of patents. Moreh, Inc. reserves the right to make corrections, modifications, enhancements, improvements, and other changes to this information, at any time and/or to discontinue any service without notice.