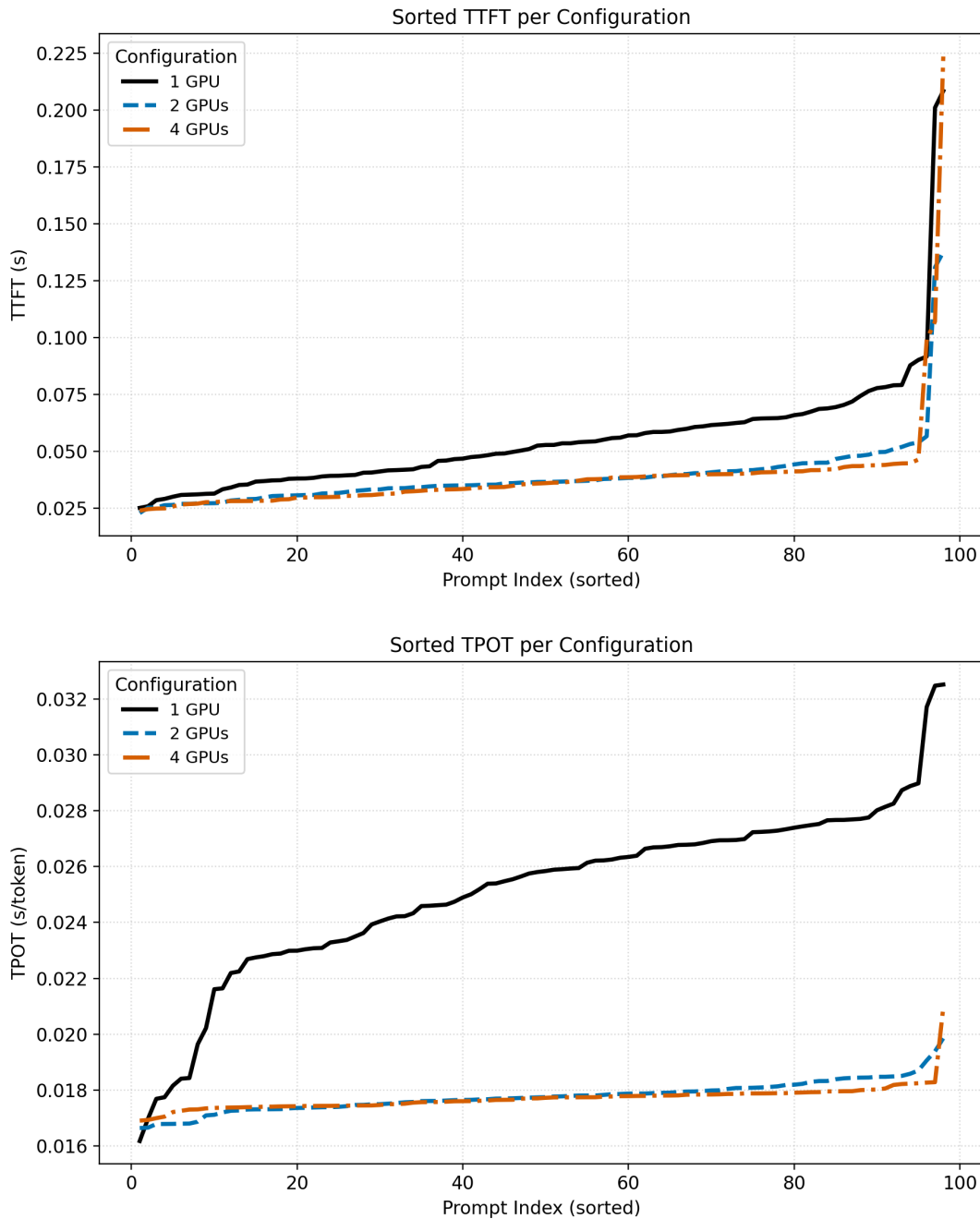


Figure of TTFT and TPOT discussion



The evaluation of TPOT (time per output token) and TTFT (time to first token) across different GPU configurations using vLLM with Llama-3.1-8B shows that a single GPU suffers from significantly higher latency, with both TPOT and TTFT increasing sharply as prompt length grows, making it impractical for real-time or long-context workloads.

In contrast, scaling to two GPUs provides a substantial improvement, flattening TPOT around 0.018–0.020 s/token and keeping TTFT consistently low, demonstrating that tensor parallelism is highly effective at this scale.

Moving from two to four GPUs yields only marginal gains, with slight stability improvements for very long prompts but limited benefits for short and medium ones, reflecting the diminishing returns caused by communication overhead.

An interesting observation from the results is that when the prompt length becomes very long, the curves for 2 GPUs and 4 GPUs cross, indicating that 2 GPUs can actually achieve lower TTFT and TPOT than 4 GPUs in this regime. This suggests that the communication and synchronization overhead introduced by scaling across more devices may outweigh the computational benefits, making 2 GPUs more efficient than 4 GPUs for extremely long prompts.

Stats System View										
CUDA API Summary	Time	Total Time	Instances	Avg	Med	Min	Max	StdDev	Name	
CUDA API Trace										
CUDA GPU Kernel Summary	24.5%	10.860 s	193700	56.063 μs	14.463 μs	4.864 μs	46.403 ms	130.960 μs	void vllm::cross_device_reduce_1stage<__nv_bfloat16, (int)2>(vllm::RankData*, vllm::RankSignals, vllm::	
CUDA GPU Kernel/Grid/Block Sum	15.6%	6.924 s	100416	68.953 μs	92.543 μs	18.847 μs	114.495 ms	35.316 μs	ampere_bf16_s16816gemm_bf16_128x64_ldg8_f2f_stages_64x4_tn	
CUDA GPU MemOps Summary (I	12.3%	5.476 s	170240	32.166 μs	25.120 μs	16.415 μs	56.640 μs	15.195 μs	ampere_bf16_s16816gemm_bf16_64x64_sliced1x2_ldg8_f2f_stages_64x6_tn	
CUDA GPU MemOps Summary (K	9.6%	4.276 s	27840	153.602 μs	153.695 μs	126.111 μs	187.486 μs	11.759 μs	void flash::flash_fwd_splitkv_kernel<Flash_fwd_kernel_traits<(int)128, (int)64, (int)128, (int)4, (bool)0, (bc	
CUDA GPU Summary (Kernels/M	7.0%	3.127 s	29236	106.970 μs	97.696 μs	94.623 μs	411.163 μs	51.816 μs	ampere_bf16_s16816gemm_bf16_256x64_ldg8_f2f_stages_64x3_tn	
CUDA GPU Trace	6.6%	2.918 s	40064	72.834 μs	56.640 μs	27.455 μs	148.159 μs	45.448 μs	ampere_bf16_s16816gemm_bf16_128x128_ldg8_f2f_stages_64x3_tn	
CUDA Kernel Launch & Exec Tim	5.9%	2.598 s	62016	41.898 μs	45.887 μs	10.368 μs	100.607 μs	25.218 μs	void flash::flash_fwd_splitkv_kernel<Flash_fwd_kernel_traits<(int)128, (int)64, (int)128, (int)4, (bool)0, (bc	
CUDA Kernel Launch & Exec Tim	2.6%	1.168 s	5376	217.286 μs	216.717 μs	30.176 μs	441.021 μs	112.366 μs	void flash::flash_fwd_splitkv_kernel<Flash_fwd_kernel_traits<(int)128, (int)64, (int)128, (int)4, (bool)0, (bc	
CUDA GPU Command List PIX R	1.8%	780.627 ms	190720	4.093 μs	4.000 μs	3.016 μs	12.608 μs	593 ns	std::enable_if_T2<(int)0&&vllm::typeConvert<T1>; exists, void=type vllm::fused_add_rms_norm_kernel<	
DX11 PIX Range Summary	1.6%	690.056 ms	150830	4.575 μs	4.192 μs	3.584 μs	355.773 μs	3.268 μs	void at::native::elementwise_kernel<(int)128, (int)4, void at::native::gpu_kernel_impl_nocast<at::native::di	
MPI Event Trace	1.3%	582.185 ms	95360	6.105 μs	6.112 μs	5.119 μs	27.263 μs	1.499 μs	void vllm::act_and_mul_kernel<(int)0:BFto16, &vllm::silu_kernel<(int)0:BFto16>, (bool)1>(T1 *, const T1	
MPI Message Size Summary	1.2%	531.223 ms	20086	26.447 μs	17.696 μs	16.639 μs	383.453 μs	54.034 μs	ampere_bf16_s16816gemm_bf16_64x64_sliced1x2_ldg8_f2f_stages_64x5_tn	
NVTX GPU Projection Summary	1.2%	524.837 ms	2032	258.286 μs	176.734 μs	112.735 μs	1.304 ms	137.993 μs	ampere_bf16_s16816gemm_bf16_256x128_ldg8_f2f_stages_64x3_tn	
NVTX Push/Pop Range Summary	1.0%	404.052 ms	99584	4.418 μs	4.768 μs	2.624 μs	6.944 μs	978 ns	void cublas::t::splitReduce_kernel<(int)32, (int)1, __nv_bfloat16, __nv_bfloat16, float, (bool)0, __nv_b	
NVTX Push/Pop Range Trace	0.9%	398.176 ms	2980	133.616 μs	130.559 μs	15.456 μs	1.735 ms	82.377 μs	ncclDevKernelAllGather_RING_LL(ncclDevKernelArgStorage<(unsigned long)4096>)	
NVTX Start/End Range Summary	0.8%	367.354 ms	95360	3.852 μs	3.808 μs	3.455 μs	10.688 μs	470 ns	void vllm::rotary_embedding_kernel<(int)0:BFto16, (bool)1>(const long *, T1 *, T1 *, const T1 *, int, long	
Network Devices Congestion	0.8%	366.709 ms	15680	23.387 μs	23.295 μs	22.688 μs	26.143 μs	426 ns	ampere_s16816gemm_bf16_128x64_ldg8_stages_64x4_tn	
MLflow API Summary	0.8%	339.160 ms	95232	3.561 μs	3.520 μs	3.072 μs	10.624 μs	427 ns	void vllm::reshape_and_cache_flash_kernel<__nv_bfloat16, __nv_bfloat16, (vllm::Fp8KVCacheDataType)0	
Users/maoxinhuan/Downloads/vllm_2gpu.suite	0.6%	268.581 ms	694	387.003 μs	388.988 μs	371.133 μs	394.876 μs	5.733 μs	ampere_bf16_s16816gemm_bf16_128x64_ldg8_f2f_stages_64x3_tn	

Identification (Screenshot) of the All-reduce kernel responsible for tensor parallel communication during multi-GPU inference

On NVIDIA GPUs, AllReduce is typically implemented by **NCCL (NVIDIA Collective Communications Library)**.

One surprising result from this assignment was the observation that **two GPUs occasionally outperformed four GPUs** in both TTFT and TPOT when handling very long prompts. Intuitively, adding more GPUs should always reduce computation time, but the profiling results revealed that the **communication and synchronization overhead** introduced at larger scales can offset these gains. Another difficulty was correctly setting up Nsight Systems to capture meaningful CUDA events, since improper configuration initially produced empty traces. Once resolved, the profiling made it clear that performance bottlenecks are not only in the compute-intensive GEMM kernels

but also in the **all-reduce communication kernel**, highlighting the challenges of balancing computation and communication efficiency in large-scale LLM inference