



Jupiter: Fast and Resource-Efficient Collaborative Inference of Generative LLMs on Edge Devices



Shengyuan Ye¹, Bei Ouyang¹, Liekang Zeng², Tianyi Qian¹, Xiaowen Chu³, Jian Tang⁴, Xu Chen¹

> ¹ Sun Yat-sen University ² The Chinese University of Hong Kong ³ The Hong Kong University of Science and Technology (GZ) ⁴ Midea Group







香港中文大學 The Chinese University of Hong Kong



Intelligent edge applications

• Large language models driven increasing intelligent applications.





Personal AI Assistants





Smart Robotics/UAV





Shengyuan Ye @ School of CSE, Sun Yat-sen University

Problems of cloud-assisted approaches

• Current LLMs-based applications heavily depend on **cloud services**.



Benefits of cloud deployment:



Powerful and scalable computing resources.

Raising three game-stopping problems:



Data privacy concerns.



Unreliable WAN connections.



Network and datacenter pressure.

Problems of on-device deployment

• On-device deployment becomes a promising paradigm for intelligent edge APPs.



 TABLE I

 INFERENCE LATENCY AND MEM. FOOTPRINT OF TRANSFORMER MODELS

WIGGET	DistilBert	Bert-L	GPT2-L	OPT-L	OPT-XL
Nano-M	0.37s	2.43s	OOM	OOM	OOM
Nvidia A100	5ms	20ms	29ms	27ms	38ms
Memory Footprint	130MB	680MB	1.6GB	2.6GB	5.4GB

121x performance gap.

Out of memory error.



Protect date privacy.



Without WAN transmission.

Limited and non-scalable onboard computing resources

Opportunities

- Edge environment often comprise a rich set of trusted idle edge devices.
 - **Opportunities**
- Smart homes usually have multiple trusted devices, such as mobile phones, laptops, and smart-home devices owned by the same family.
- ✓ Utilize nearby edge devices as resource augmentation to render expedited LLMs inference at the edge.





Smart Phones/Watches



Personal Laptops/Server



Edge Computing Platforms

Distributed Inference Methods for LLMs

 DP, SP, TP, and PP represent the four most prominent techniques for parallel inference in large language models (LLMs).



- The primary goal of edge inference systems is to reduce the latency of singlesequence requests from users. <u>Common DP or PP fail to meet the requirements</u>.
- Existing collaborative edge inference systems predominantly use TP or SP as the major parallel architecture, as exemplified by Galaxy [1] and DeTransformer (DT) [2].



[1] Ye et al., "Galaxy: A Resource-Efficient Collaborative Edge AI System for In-situ Transformer Inference". IEEE International Conference on Computer Communications (INFOCOM), 2024.

[2] Wei et al., "Communication-Efficient Model Parallelism for Distributed In-situ Transformer Inference". Design, Automation and Test in Europe Conference (DATE), 2024.

- TP/SP-based architectures require **multiple collective communications** within each decoder layer, leading to significant communication overhead.
- In each decoder layer, TP requires two AllReduce operations to synchronize tensors, while SP requires two AllGather operations to collect keys and values.



An instance of tensor model parallelism across two edge devices

Shengyuan Ye @ School of CSE, Sun Yat-sen University

• TP/SP-based architectures require **multiple collective communications** within each decoder layer, leading to significant communication overhead.

TABLE II

COMM.-TO-COMP. RATIO OF VARIOUS PARALLELISM METHODS.

Model	Network	Communication-to-Computation Ratio							
Name	Bandwidth	SP	TP	DT [7]	Galaxy [6]	Jupiter			
Llama2-7B	100Mbps	8.16	6.96	3.48	5.19	0.08			
	1Gbps	0.92	0.88	0.45	0.69	0.01			
Llama2-13B	100Mbps	5.71	6.06	3.03	4.63	0.05			
	1Gbps	0.73	0.81	0.38	0.56	0.01			

Existing TP/SP-based systems exhibit a high communication-tocomputation latency ratios making communication a critical bottleneck.

Up to 92.6x

Architecture Choices of Jupiter

• Jupiter employing a **pipeline architecture** to orchestrate multiple edge devices.



Observation: Pipeline architecture has a *lower communication-tocomputation ratio* during collaborative LLM inference, making it more suitable for low-bandwidth edge environments.



Architecture Choices of Jupiter

• Jupiter employing a **pipeline architecture** to orchestrate multiple edge devices.

Challenges: Pipeline architecture struggles to accelerate inference for single-sequence requests *due to the lack of batch dimension*.



- Existing research works are mostly limited to discriminative tasks and focus solely on the prefill phase.
- The main strength of modern LLMs lies in generative decoding tasks, but most existing systems are **not designed or optimized for the decode phase**.



Challenges

• Efficiently leveraging the computational resources of multiple edge devices under a **pipeline architecture** is a non-trivial problem.



1. How to accelerate the **prefill phase** with pipelined architecture for single sequence request.

2. How to accelerate the autoregressive decoding phase with pipelined architecture.



Challenges

• Efficiently leveraging the computational resources of multiple edge devices under a **pipeline architecture** is a non-trivial problem.



1. How to accelerate the **prefill phase** with pipelined architecture for single sequence request.

2. How to accelerate the **autoregressive decoding phase** with pipelined architecture.



Collaborative Inference for Prefilling Phase

• Solution to Challenge #1: Utilizing Intra-sequence pipeline parallelism.

- 1. Partition the LLM into multiple stages, each handled by a separate edge device.
- 2. Partition the input sequence into multiple sub-sequences to increase parallelism.
- 3. When computing the sub-sequence s_i , the **cached key and value** of s_1, \ldots, s_{i-1} are utilized to ensure accurate self-attention results.



Fig. 5. An illustration of pipelined inference with three edge devices.

Fig. 6. An illustration of opportunities of intra-sequence parallel inference.

Optimal Model and Sequence Partition

Selecting Optimal LLMs Partition

• Use dynamic programming algorithm to partition the target LLM into multiple stages with **equal execution time**.



Optimal Model and Sequence Partition

• Selecting Optimal Sequence Partition

Sequence Partition is not trivial for the following reasons:

- Increasing the number of sub-sequences can efficiently boost parallelism. However, this results in shorter sub-sequences, which may underutilize the mobile accelerators.
- 2. Partitioning the sequence into sub-sequences of equal length for pipelining is not optimal. Later sub-sequences carry a heavier computational load than earlier ones.
- 3. The input sequence **lengths vary across different requests**, necessitating the determination of the optimal partitioning strategy for any possible length prior to online service deployment

Optimal Model and Sequence Partition

Selecting Optimal Sequence Partition

 We carefully profile and model the inference latency and design a dynamic programming algorithm to determine the optimal partitioning for requests with varying sequence lengths.



Fig. 5. An illustration of pipelined inference with three edge devices.

Collaborative Inference for Decoding Phase

• Solution to Challenge #2: Integrate speculative decoding with pipeline architecture



Observation: Compared to speculative decoding with small-large model collaboration, self-speculative decoding is better suited for distributed pipelined deployment.



Fig. 8. A workflow of our collaborative inference with speculative decoding.

Collaborative Inference for Decoding Phase

Solution to Challenge #2: Outline-Based Pipeline Parallel Decoding

Observation:

- LLMs tend to generate structured, point-by-point responses.
- The generation of each individual point is often *independent* and lacks contextual coherence.



Fig. 9. An illustration of our outline-based pipeline parallel decoding.

Putting It All Together

Jupiter system workflow.

Jupiter uses a pipelined architecture to significantly reduce communication overhead. It optimizes not only the prefill phase, but also the decoding phase.



• Testbeds

Using these 3 heterogeneous devices, we simulated 2 different edge clusters, including both homogeneous and heterogeneous clusters.

TABLE III SPECIFICATIONS OF EDGE DEVICES IN EXPERIMENTS.

Edge Device	GPU Processor	Memory	Power
Jetson Xavier NX [16]	384-core NVIDIA Volta	8GB	20W
Jetson TX2 [30]	256-core NVIDIA Pascal	8GB	20W
Jetson Nano [31]	128-core NVIDIA Maxwell	8GB	10W

Env. A (Homo.): 4 x NX

Env. B (Hetero.): 1 x NX + 2 x TX2 + 1 x Nano

Models and datasets

- 2 LLMs from the Llama2 series, specifically Llama2-7B and Llama2-13B (both with INT4 quantization)
- Evaluate with 3 recent assistant-style datasets: LiMA, Vicuna-80, and WizardLM.

Baselines

- > Megatron-LM (M-LM): A state-of-the-art tensor model parallelism method.
- Sequence Parallelism (SP): A state-of-the-art sequence model parallelism method
- DeTransformer (DT): TP-based collaborative edge system that optimizes communications by reducing the frequency of tensor synchronization.
- Galaxy: TP-based collaborative edge system that optimizes communications by fine-grained overlapping of comm. and comp.
- EdgeShard (ES): PP-based collaborative edge inference system that employs pipelined architecture to orchestrate edge devices.



64 tokens

Jupiter achieves up to **26.1x end-to-end generation latency reduction** compared to baselined system!

	Edge	Network				Llama2-7	В				L	lama2-13	В	
Average input	Environment	Bandwidth	SP	M-LM	DT	Galaxy	EdgeShard	Jupiter	SP	M-LM	DT	Galaxy	EdgeShard	Jupiter
sequence length:		100Mbps	53.5	431.2	228.5	427.6	42.2	16.5	OOM	503.4	270.1	496.5	66.2	26.3
260 tokensHomo.Env. A	Homo. Env A	500Mbps	37.4	106.9	66.4	103.9	39.0	15.2	OOM	130.1	83.4	125.0	63.4	25.2
	Liiv. A	1Gbps	35.4	66.4	46.1	65.0	38.6	14.9	OOM	83.4	60.1	81.3	63.1	24.9
Maximum		100Mbps	63.1	491.2	288.6	458.3	59.3	22.4	OOM	624.5	391.2	566.4	102.4	38.8
	Hetero. Env B	500Mbps	47.0	167.0	126.4	142.9	56.1	21.4	OOM	251.2	204.5	208.0	99.7	37.3
generation length:	h:	1Gbps	44.8	126.4	106.2	104.9	55.7	20.9	OOM	204.5	181.2	165.7	98.3	36.8





Fig. 10. Evaluate in Homogeneous Environment A. The average per-token processing/generation latency in prefill/decoding phase. \times indicates OOM.

Fig. 11. Evaluate in Heterogeneous Environment B. The average per-token processing/generation latency in prefill/decoding phase. \times indicates OOM.



The high comm.-to-comp. ratio of baselines hinders resource-efficient scaling in bandwidth-limited edge environments. Jupiter outperforms baseline methods in **scalability** across varying bandwidths.





Jupiter achieves up to **3.9x speedup** over naive sequential generation while maintaining comparable generation quality.

TABLE V
Speedup over naive sequential generation. SD: Speculative
DECODING. OP: OUTLINE-BASED PARALLEL DECODING.

Madal	Speedup Over Naive							
WIGUEI	Naive	Jupiter w/o OP	Jupiter w/o SD	Jupiter				
Llama2-7B	$1.0 \times$	$1.8 \times$	2.3 imes	$3.6 \times$				
Llama2-13B	$1.0 \times$	2.0 imes	2.4 imes	$3.9 \times$				

TABLE VIOVERALL ANSWERS QUALITY OF NAIVE AND JUPITER'S METHOD.

Mothod-	Vicu	na-80	Wiza	rdLM	LiMA		
Wiethou	Llama2-7B	Llama2-13B	Llama2-7B	Llama2-13B	Llama2-7B	Llama2-13B	
Naive	7.05	7.19	6.26	6.85	6.72	6.77	
Jupiter	6.59	6.85	6.21	6.70	6.20	6.25	

Jupiter Demo



Run collaborative inference of the Llama2-7B model using four NVIDIA Jetson edge devices.









Thanks for listening

Shengyuan Ye¹, Bei Ouyang¹, Liekang Zeng², Tianyi Qian¹, Xiaowen Chu³, Jian Tang⁴, Xu Chen¹

¹ Sun Yat-sen University
 ² The Chinese University of Hong Kong
 ³ The Hong Kong University of Science and Technology (GZ)
 ⁴ Midea Group





香港中文大學 The Chinese University of Hong Kong

