

Early Experience in Characterizing Training Large Language Models on Modern HPC Clusters

Hao Qi, Liuyao Dai, Weicong Chen, and Xiaoyi Lu

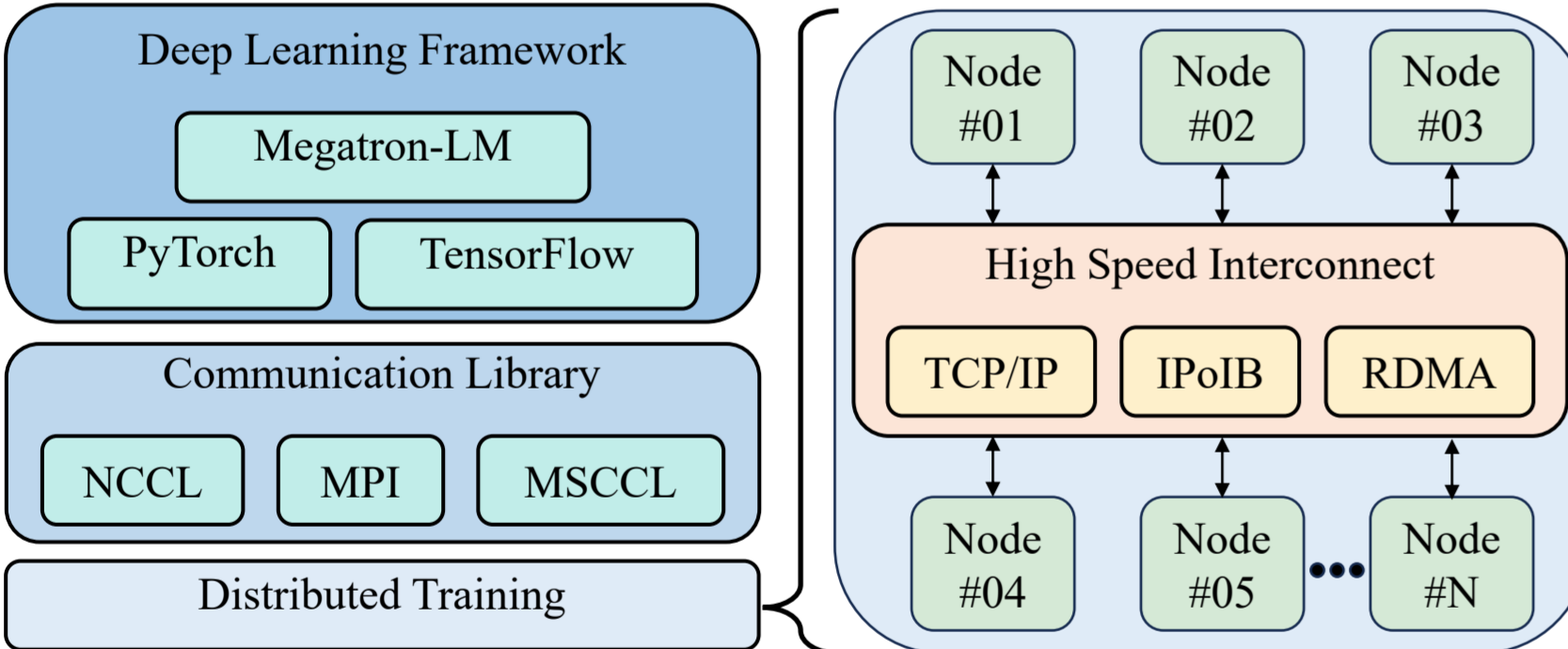
University of California, Merced



ABSTRACT

- Large Language Models (LLMs) are powerful tools in natural language processing, aiding in tasks like language translation, text generation, and sentiment analysis.
- The immense size and complexity of LLMs pose significant challenges.
- This study explores high-performance interconnects in the distributed training of various LLMs.
- High-performance network protocols, especially RDMA, are found to significantly outperform IPoIB and TCP/IP in training performance.
 - RDMA offers improvements by factors of 2.51x and 4.79x respectively.
- Despite the significant findings, there is potential for further optimization in overall interconnect utilization.
- The research provides deeper insights into the performance characteristics of LLMs over high-speed interconnects and paves the way for developing more efficient training methodologies for LLMs.

INTRODUCTION



Challenges in Distributed Training for LLMs:

- Need for communication & coordination among nodes and GPUs.
- Vast amount of training data.
- Requirement for GPU-enabled distributed training.
- Emphasized need for high-performance interconnects.

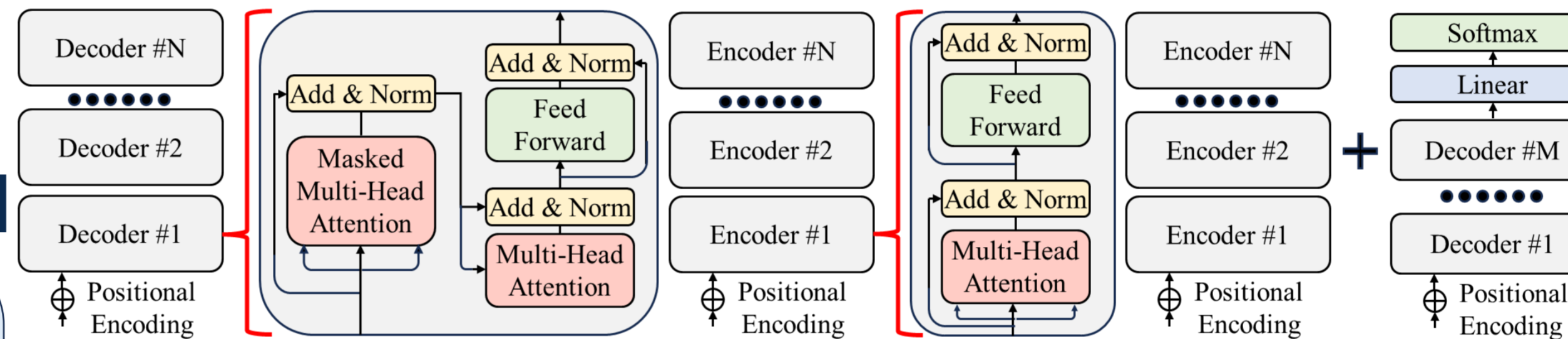
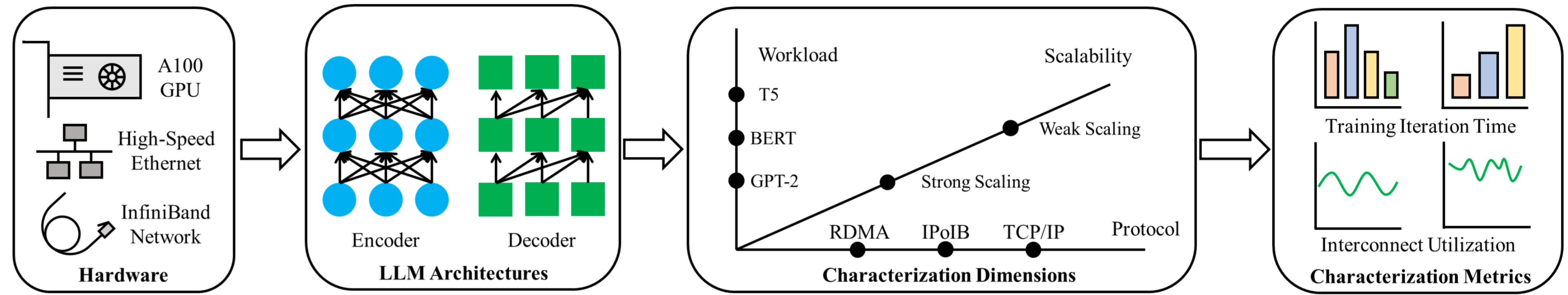
Role of High-speed Interconnects:

- Ensure efficient data transfer and synchronization.
- Crucial for rapid and scalable communication among nodes and GPUs.

Research Questions:

- Potential bottlenecks in communication and the time proportion for communication in various LLM training configurations.
- Efficiency of current high-performance interconnects in different distributed training scenarios.
- Quantitative performance impact of networking technologies & protocols (e.g., RDMA, IPoIB, TCP/IP) on LLM training.

METHODOLOGY



Characterization Dimensions: Workload, Interconnect/Protocol, and Scalability.

- Workload:** GPT-2-Medium, GPT-2-Large, BERT-Large, and T5-Large.
- Interconnect/Protocol:** TCP/IP, IPoIB, and RDMA (with GPUDirect).
- Scalability:** Strong scaling and Weak scaling, and maximum batch size.
- Framework:** Megatron-LM.
- Dataset:** enwiki dataset (20.4 GB).

EVALUATION

Fraction of Comm. vs. Iteration Time	Workload	Protocol	Strong Scaling			Weak Scaling			Larger Batch Size		
			<2,16>	<4,16>	<8,16>	<2,8>	<4,16>	<8,32>	<8,96>	<8,256>	<8,512>
T5-Large	TCP/IP		0.25	0.83	0.9	0.27	0.83	0.9	0.87	OOM	OOM
	IPoIB		0.24	0.74	0.78	0.27	0.74	0.78	0.74	OOM	OOM
	RDMA		0.24	0.42	0.47	0.27	0.42	0.46	0.41	OOM	OOM
BERT-Large	TCP/IP		0.23	0.85	0.91	0.3	0.85	0.91	0.84	0.7	OOM
	IPoIB		0.23	0.76	0.81	0.3	0.76	0.8	0.68	0.49	OOM
	RDMA		0.23	0.46	0.52	0.3	0.46	0.49	0.34	0.3	OOM
GPT-2-Large	TCP/IP		0.27	0.88	0.94	0.35	0.88	0.93	0.86	OOM	OOM
	IPoIB		0.27	0.81	0.87	0.36	0.81	0.84	0.72	OOM	OOM
	RDMA		0.27	0.53	0.61	0.35	0.53	0.57	0.38	OOM	OOM
GPT-2-Medium	TCP/IP		0.24	0.86	0.92	0.3	0.86	0.91	0.85	OOM	OOM
	IPoIB		0.24	0.78	0.84	0.31	0.78	0.81	0.69	OOM	OOM
	RDMA		0.24	0.48	0.54	0.31	0.48	0.5	0.35	OOM	OOM

RDMA-100 Gbps outperforms IPoIB-100 Gbps and TCP/IP-10 Gbps by an average of 2.51x and 4.79x regarding training iteration time, and scores the highest interconnect utilization (up to 60 Gbps) in both strong and weak scaling, compared to IPoIB with up to 20 Gbps and TCP/IP with up to 9 Gbps, leading to the shortest training time.

CONCLUSION

Main Contribution:

- Enhanced understanding of performance dynamics of LLMs over high-speed interconnects.

Research Efforts:

- Detailed exploration of communication's role in distributed LLM training.
- Findings can guide the design and setup of systems for growing LLM application demands.

Evaluation Outcomes:

- Both strong scaling and weak scaling experiments show similar patterns.
- Highlighted impact of interconnect/protocol on distributed training.
- Emphasized the importance of efficient interconnect utilization.

Directions for Future Work:

- Explore alternative parallelism strategies, e.g., model parallelism.
- Study behavior of even larger models at greater scales.
- Develop methods to further enhance interconnect utilization.

ACKNOWLEDGEMENT

This work was supported in part by an NSF research grant OAC #2321123, a DOE research grant DE-SC0024207, and an Amazon Research Award. Part of this research was conducted using Pinnacles (NSF MRI, #2019144) at the Cyberinfrastructure and Research Technologies (CIRT) at the University of California, Merced.