## SymTensor: Symbolic and Adaptive Tensor Partitioning by Unified Parallelism for Deep Learning

Hongxing Wang, Zhengdao Yu, <u>Chong Li</u>, Serge Petiton

HLPP 2025 @ Innsbruck



# Background: Rapid Evolution of Large Language Models (LLM)



Image source:Harnessing the Power of LLMs in Practice: A Survey on ChatGPT and Beyond, J JIANG et al., AMAZON USA, 2023

LLM evolution remains extraordinarily rapid:

- Expansion of parameters: - From millions to hundreds of billions

### - Structural diversity:

- BERT -> Encoder-Decoder architecture
- GPT -> Decoder-only architecture
- Mixtral -> MoE architecture
- DeepSeek -> Dense + MoE architecture and more LLM...

### - Operator innovations:

- Self-Attention -> FlashAttention
- Operator Fusion: Mul + Add -> Mul\_Add ...

LLM designers want a software framework that could:

- 1. help them to focus on DL design
- 2. simplifying performance tuning (on distributed clusters)

# Deep Learning (DL) Frameworks: Designed for User-Friendly

## 1. Simplifying for Training



## 2. Simplify for Parallelization

DL Framework	Parallelization Planner		
Pytorch	Megatron-LM		
Jax	Alpa		
MindSpore	SAPP		

Planner: decide & manage Comm. Operator

designer's perspective:



Indirect control necessitates systematic approaches

# Parallelisms in DL: Predefined Patterns with Performance Guarantee

#### • Data Parallelism (DP):

Replicate model on each worker, split the input data



### • Tensor Parallelism (TP):

Split weight of operators (MatMul) across workers



### Skip bad choices, e.g. for X \* Y \* Z:



Hide comm. cost:

## Strategy of Parallelisms: Finding Best Performance for End-User

	Definition	Example
Strategy	Configure degrees of parallelism over the AI model	(DP = 4, TP = 2) over 8 devices

(DP = 4, TP = 2) over 8 devices:



### Experiment on Deepseek V3 671B (64 devices)

MFU - Model FLOPs Utilization

Config	DP	ТР	MFU	Insights
Option 1	8	8	26.84%	Balanced DP/TP
Option 2	64	1	Out-of- Memory	DP-only parallelism $\rightarrow$ OOM
Option 3	32	2	31.42%	High-DP within memory limits -> Better performance

### Parallelism Strategy is important:

- Same 64 devices, different DP/TP configs -> vastly different outcomes
- From OOM to 26.84% to 31.42% MFU -> strategy choice makes the difference

# Challenges

## **Challenge 1: Existing frameworks lack adaptability**

## for evolving models & parallelisms

Relying on predefined rules and paradigms: DP, TP, ...

### **Problem:**

New parallelisms continue to emerge to meet evolving training demands

- Sequence Parallelism (SP) for long sequence scenario (4K tokens -> 128K tokens)
- Expert Parallelism (EP) for MoE architecture



## **Challenge 2: Difficult to determine best parallelism**

### strategy

#### **Parallelisms were designed independently** However parallelism strategies require combining them

### Problem:

Parallelism combinations yield undetermined performance impacts

### Experiment on Deepseek V3 671B (64 devices)

Config	DP	ТР	MFU	Insights
Option 1	8	8	26.84%	Balanced DP/TP
Option 2	64	1	Out-of- Memory	DP-only parallelism $\rightarrow$ OOM
Option 3	32	2	31.42%	High-DP within memory limits -> Better performance

Lack of comprehensive cost model

## **Our Solution**

**For challenge 1** *Lack of unified abstraction of parallelisms* 

## **For challenge 2** Lack of comprehensive cost model



# Solution in Detail



# **Experiment Environment**

**Open-Source DL framework:** https://gitee.com/mindspore/mindspore

## MindSpore



### Ascend 910 Cluster

2048 Node x 256TFlops = 512 Peta Flops



Scalability with Hierarchical Comm has been discussed in: H. Wang, C. Li, T. Tachon, et al., "Efficient and systematic partitioning of large and deep neural networks for parallelization," in Euro-Par 2021

This paper focuses on Comm- & Mem-aware

Exp str on a 910 A2 node

8x Ascend NPU per node Each NPU:

- To CPU/Mem: PCle 4.0 x16
- To neighbors: -

Grouped AllGather

 $\begin{bmatrix} X_{01} \\ X_{11} \end{bmatrix}$ 

 $\begin{bmatrix} X_{01} \\ X_{11} \end{bmatrix}$ 

- 7x 56 GB/s full mesh HCCS
- 1x 200G Eth



# **Cost Model Validation**



#### **Best choice of strategy (Memory)**

14.38

AllReduce

14.3

## Memory

- Memory-aware cost model ٠
  - Precisely predicts the lowest relative cost strategy ("\*")



## Communication

- SymTensor captures the relative costs between • different strategies
  - Lowest predicted cost matches the optimal strategy choice

SymTensor focuses on choices -> tolerated in a good level of precision

**Best choice of strategy (Communication)** 

# **Real-case Validation**

## Generality

#### **Experiment objective:**

- Evaluate SymTensor's effectiveness

### Experimental setup:

we chose 6 typical DL models for:

- Dense autoregressive language models: LLaMA3-8B, LLaMA3-70B
- Instruction-tuned model for code generation: Deepseek-Coder-7B
- LLM with standard transformer backbones: InternLM2-20B, CodeLLaMA-34B
- Bilingual Model using Lora:

Qwen 7b

#### Speedup of SymTensor vs. Baseline (Megatron-LM: Optimized and Tuned Strategy)



## Adaptability

#### Experiment objective:

- Demonstrate adaptability to common real-world model variations

#### Experimental setup:

we chose widely used foundation models for:

- Operator substitution (LLama2-13B): Replace self-attention -> FlashAttention
- New MoE architecture (Mixtral 8x7B): Novel MoE design
- High Memory scenario (Qwen-7B-LoRA): Increase the batch size from 8 to 32

### Table: Training Throughput Comparison (in tokens/sec)

Model	Megatron- LM	SymTensor	SpeedUp
LLaMA2 13B	10961	15845	144.56%
Mixtral <sup>8x7B</sup>	2936	6506	221.59%
<b>Qwen</b> 7B LoRA	OOM	17626	∞

# Conclusion & More...

We demonstrated how to systematically optimize:

- Training -
- for LLMs \_
- over typical AI-accelerator Clusters -

We are open to further research challenges on :

- New architectures (superPod, ...) -
- New AI model (multi-modal, ...) -
- Other distributed DL scenario (inference, ...) -
- And all topics related to parallel and AI (for perf, ...) -

Heterogeneity-aware Placement

Manager

Submodule

partitions

Data

Bs:3×24

Bs: 4x18







https://arxiv.org/pdf/2506.12708

## Huawei CloudMatrix384

