# **Mini-Sequence Transformer:** Optimizing Intermediate Memory for Long Sequences Training

Cheng Luo[1], Jiawei Zhao[2], Zhuoming Chen[3], Beidi Chen[3], Anima Anandkumar [1]
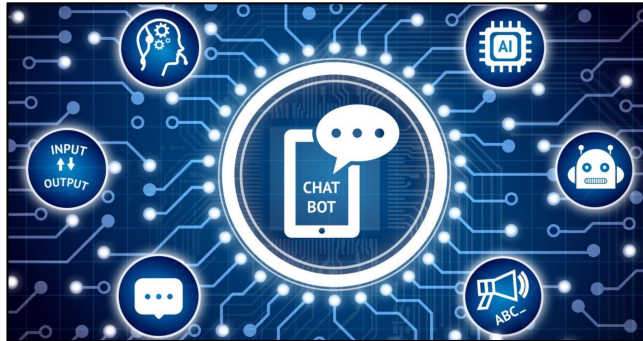
California Institute of Technology[1], Meta Fair[2], Carnegie Mellon University[3]

wdlctc@gmail.com
https://github.com/wdlctc/mini-s

# Foundation Models



Conversational AI



Content Generation



AI Agents



Reasoning



Planning

# Evolution of Foundation Models
# LLama 3 vs. LLaMa 2 - A Head-to-Head Comparison

| Feature | LLAMA2-7b | LLAMA3-8b |
|---|---|---|
| Hidden Length | 4096 | 8192 |
| Sequence Length | 4096 | **8192** |
| MLP Intermediate | 11008 | **14336** |
| Vocal Length | 32000 | **128256** |

**Evaluation of foundataion model leads to larger Vocal Length & MLP Intermediate**

# Mini-Sequence Transformer Yield Memory Saving & Long Sequence Enable

Loss

LM head

L x

MLP

Norm

Attention

Norm

Tensor (S)

MLP/LM-head Computation Flow
Extremely Large Intermediate

MLP/LM_head

X → Intermediate → O

**Challenge:
Native implementation
brought memory bottleneck**

Llama-3 Sandard Implementation:
Peak Memory:67GB

MLP/LM_head  - MsT

X → Intermediate → O

MsT: Splitting MLP/LM-Head Input
        Concat output

Llama-3 MST Implementation:
Peak Memory:47GB
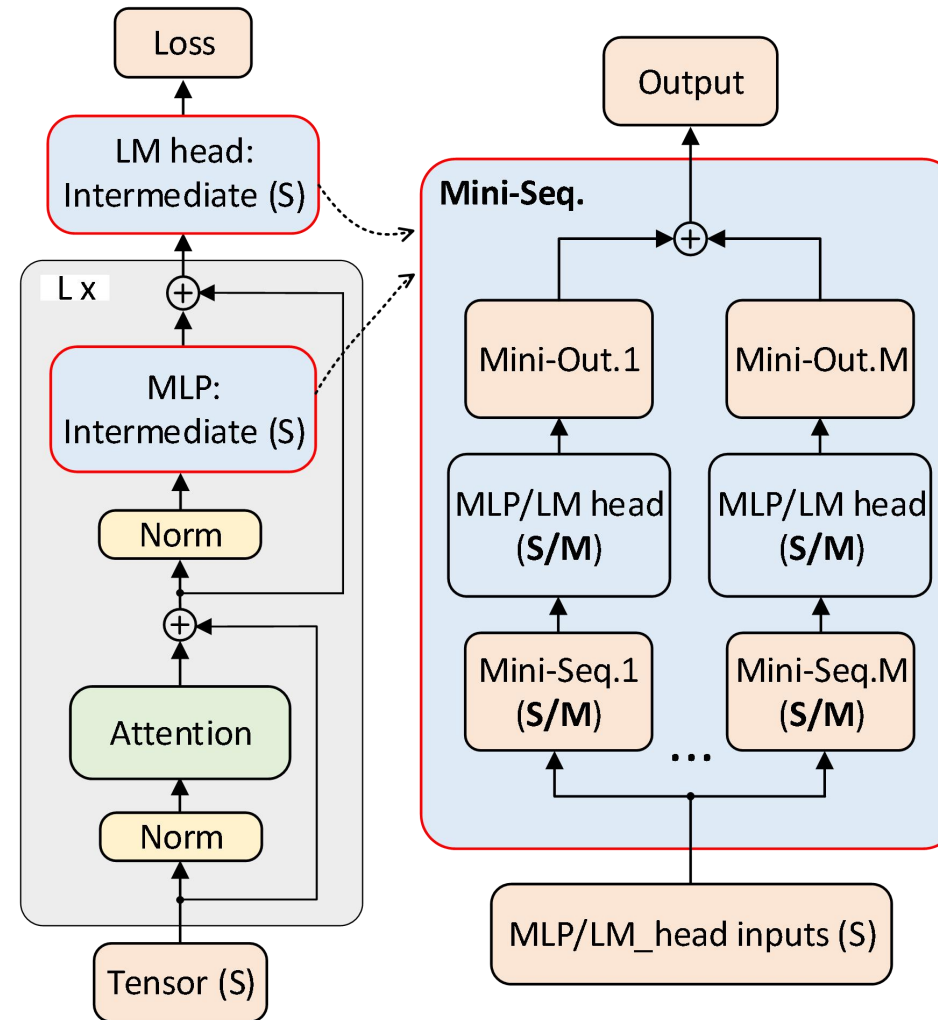
**Gain: Support 12x-24
longer Sequence Length**

# Mini-Sequence Transformer: Key Ideas

Key Ideas: Input Partition & Gradient Accumulation

Input Partition: Split the input sequence into Mini-sequence

Gradient Accumulation: Restore the full output/Gradients

Results: optimizing intermediate values, create more space for long sequence activation



**Transformer -> Mini Sequence Transformer**

# Mini-Sequence Transformer (MsT): 12-24x sequence & TFLOPS equivalence

12-24x longer sequence enable

| Models | Maximum Sequence (K) | | |
|---|---|---|---|
| | Vanilla | MST | Extend |
| LLAMA3-8b | 5 | **60** | **12x** |
| LLAMA2-78 | 7 | **84** | **12x** |
| Qwen2-7B | 4 | **74** | **18x** |
| gemma-2-9b | 1.5 | **36** | **24x** |

Equal TFLOPS for Llama3-8B

| Models | TFLOPS | | |
|---|---|---|---|
| | Vanilla | MST | Speedup |
| LLAMA3-8b | 3271 | **3386** | **1.02x** |
| LLAMA2-78 | 3290 | **3656** | **1.11x** |