# MEGATRON

Large-Scale Language Model Training Framework

## TRAINING FRAMEWORK [megatron/training/]

| Training Loop | Arguments | Checkpointing | Initialization |
|---|---|---|---|
| *training.py* | *arguments.py* | *checkpointing.py* | *initialize.py* |

## MODELS `MCore`

| GPT | BERT | T5 | Mamba | Multimodal | MoE |
|---|---|---|---|---|---|
| *core/models/gpt/* | *core/models/bert/* | *core/models/t5/* | *core/models/mamba/* | *core/models/multimodal/* | *core/models/mimo/* |

## TRANSFORMER CORE

| Attention | MLP | Layer Norm | Embeddings | Block |
|---|---|---|---|---|
| *core/transformer/attention.py* `MCore` | *core/transformer/mlp.py* `MCore` | *core/transformer/* `MCore` | *core/transformer/* `MCore` | *core/transformer/* `MCore` |

## PARALLELISM STRATEGY [megatron/core/]

### TENSOR PARALLELISM
Intra-layer model sharding

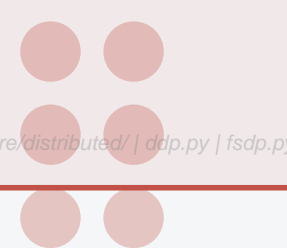*core/tensor_parallel/ | layers.py | mappings.py*

### PIPELINE PARALLELISM
Inter-layer distribution

### DATA PARALLELISM
Gradient synchronization

*core/distributed/ | ddp.py | fsdp.py*

## DATA PIPELINE

| Blended Dataset | Indexed Dataset | GPT Dataset | Data Loader |
|---|---|---|---|
| *core/datasets/blended_* | *core/datasets/indexed_* | *core/datasets/gpt_* | *training/datasets/* |

## OPTIMIZER

| Distributed Optimizer | Adam Config | Gradient Scaling |
|---|---|---|
| *core/optimizer/* `MCore` | *core/optimizer/* `MCore` | *core/optimizer/* `MCore` |

## CHECKPOINTING [training/]

| Save/Load | Distributed | Model State |
|---|---|---|
| *training/* | *core/dist_checkpointing/* `MCore` | *training/* |

## KEY FEATURES [megatron/core/]

### PRECISION SUPPORT `MCore`
- FP16 / BF16 Training
- FP8 (Hopper Optimized)
- FP4 Quantization
- Transformer Engine

### MEMORY OPTIMIZATION `MCore`
- Activation Checkpointing
- Sequence Parallelism
- Gradient Checkpointing
- Fine-grained Offloading

### PERFORMANCE `MCore`
- CUDA Graphs
- Fused Kernels
- Flash Attention
- Fused LayerNorm

## ARCHITECTURAL PATTERNS

### Modular Core Library `MCore`
megatron/core/ provides production-ready,
GPU-optimized building blocks for
framework developers

*models | transformer | tensor_parallel*
*pipeline_parallel | distributed | optimizer*

### Configuration-Driven
Dataclass-based configuration system
enables flexible model parallelism
setup via command-line arguments

*TransformerConfig | ModelParallelConfig*
*ProcessGroupCollection | parallel_state*