# Distributed GPT2 Model Project

Abhijeet Sahdev, Andrew Mathews, Felix Guerrero

# Rationale

- AI is the biggest topic in computing over the last few years

- GPU parallelism allows thousands of operations to run at once

- Training AI can take a really long time on a single machine

- Using MPI and/or CUDA can dramatically reduce the time required for training an AI model

- This project would help us understand how to train an AI model

# Methodology

➔ We briefly went over:
- ◆ Pytorch whitepaper and documentation
  - ● [https://arxiv.org/pdf/1912.01703](https://arxiv.org/pdf/1912.01703))
  - ● Simple to use deep learning framework
- ◆ Build a Large Language Model (From Scratch) by Sebastian Raschka
- ◆ Let's build GPT: from scratch, in code, spelled out. A youtube video by Andrej Karpathy

➔ Our C implementation is heavily inspired by Pytorch object-oriented approach

# Codebase

- Every stage has the following files:

| | | |
|---|---|---|
| Makefile | dataloader.h | run_inference.sh |
| README.md | embedding.h | softmax.h |
| adam.h | gelu.h | tensor.h |
| add.h | gpt.h | tokenizer.h |
| autograd.h | head.h | train_gpt2.cpp |
| block.h | inference.cpp | trained_weights.bin |
| broadcast.h | layernorm.h | transpose.h |
| checkpoint.h | linear.h | |
| cross_entropy.h | matmul.h | |
| cuda_kernels.cu | mlp.h | |
| cuda_utils.h | multihead_attention.h | |

# Differences Between Versions

- **Serial**: Embeddings, transformer blocks, output head
- **MPI**: Batch sharing, parameter sync metadata, distributed generation, gpt_set_distributed
- **CUDA**: Core operations offloaded to kernels, no tokenizer/dataloader broadcast, same model graph
- **MPI + CUDA**: Rank 0 builds tokenizer and broadcasts parameters, coordinated distributed settings across processes and cores, end-to-end flow similar to CUDA version, but with mandatory data-parallel synchronization

# Figures

## Model Architecture

Input token IDs $(B, T)$

Token embedding
wte : vocab_size $\times n_{embd}$

Position embedding
wpe : $T \to n_{embd}$

broadcast_add $\to x_0 \in R^{B \times T \times n_{embd}}$

Transformer block $\times n_{layer}$

LayerNorm 1
$x \mapsto \mathrm{LN}_1(x)$

Multi-Head Attention
Per head: $q = W_q x$, $k = W_k x$, $v = W_v x$
$A = \dfrac{qk^\top}{\sqrt{\text{head\_size}}}$, causal mask $t_k > t_q \Rightarrow -10^{30}$
$\alpha = \mathrm{softmax}(A)$, head_out $= \alpha v$
concat heads $\to$ Linear proj $\to$ attn_out

Residual:
$x + $ attn_out

LayerNorm 2
$x_1 \mapsto \mathrm{LN}_2(x_1)$

Residual:
$x_1 + $ mlp_out

MLP
$c_{fc} : n_{embd} \to 4n_{embd}$
GELU
$c_{proj} : 4n_{embd} \to n_{embd}$
$\to$ mlp_out

Next $x$ (block output)

Final LayerNorm
$x \mapsto \mathrm{LN}_f(x)$

lm_head : $n_{embd} \to$ vocab_size
$\Rightarrow$ logits $\in R^{B \times T \times V}$

## Stage 0

data/dummy.txt

tokenizer_extract + tokenizer_encode
+ tokenizer_pad_to
*Tokenizer (vocab, encoded ids, PAD/EOS)*

DataLoader.next_batch
$\to$ inputs[B*T], targets[B*T]

gpt_forward_with_loss (GPT model)
$\to$ logits (B,T,V)
$\to$ cross_entropy_loss_3d $\to$ loss (1)

backward(loss)
*autograd over tracked tensors*

adam_step
*clip_grad_norm, lr schedule on collected params*

adam_zero_grad

save_weights("trained_weights.bin")

generate_sample_text
$\to$ gpt_forward_logits
$\to$ greedy_select_next_token

**loop over epochs**
+ batches

## Stage 1

**MPI init & RNG**
MPI_Init, parse --seed
per-rank RNG via tensor.set_seed(seed + rank)

**Tokenizer on rank 0**
Tokenizer("../data/dummy.txt")
extract $\to$ encode $\to$ pad

**Broadcast tokenizer**
tokenizer_broadcast
all ranks reconstruct vocab/ids

**DataLoader init**
global batch $B$, per-rank $B/world\_size$
init_with_tokenizer

**Broadcast dataloader**
shared sampling state

**Model init + distributed config**
gpt_init, gpt_set_distributed
collect params $\to$ ParamSyncInfo

**Initial parameter sync**
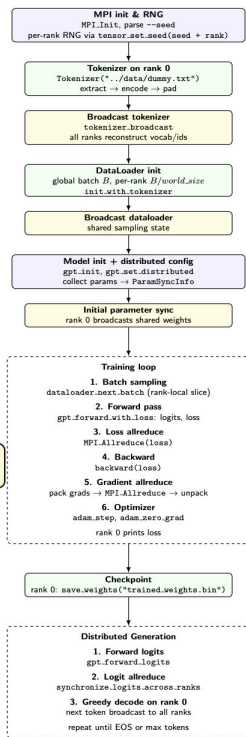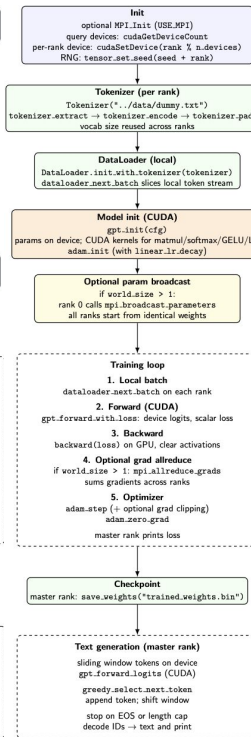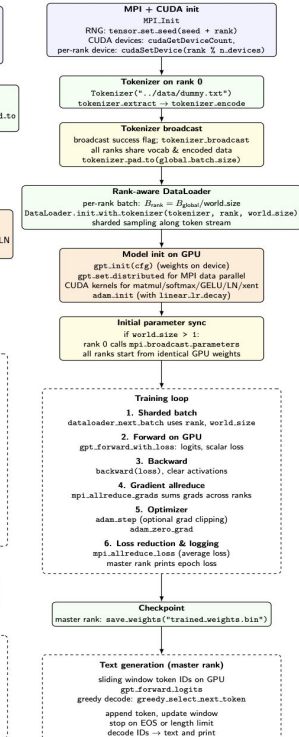rank 0 broadcasts shared params

**Training loop**
**1. Batch sampling**
dataloader.next_batch (rank-local slice)
**2. Forward pass**
gpt_forward_with_loss: logits, loss
**3. Loss allreduce**
MPI_Allreduce(loss)
**4. Backward**
backward(loss)
**5. Gradient allreduce**
pack grads $\to$ MPI-Allreduce $\to$ unpack
**6. Optimizer**
adam_step, adam_zero_grad
rank 0 prints loss

**Checkpoint**
rank 0: save_weights("trained_weights.bin")

**Distributed Generation**
**1. Forward logits**
gpt_forward_logits
**2. Logit allreduce**
synchronize_logits_across_ranks
**3. Greedy decode on rank 0**
next token broadcast to all ranks
repeat until EOS or max tokens

## Stage 2

**Init**
optional MPI_Init (USE_MPI)
query devices: cudaGetDeviceCount
per-rank device: cudaSetDevice(rank % n.devices)
RNG: tensor.set_seed(seed + rank)

**Tokenizer (per rank)**
Tokenizer("../data/dummy.txt")
tokenizer_extract $\to$ tokenizer_encode $\to$ tokenizer_pad_to
vocab size reused across ranks

**DataLoader (local)**
DataLoader.init_with_tokenizer(tokenizer)
dataloader.next_batch slices local token stream

**Model init (CUDA)**
gpt_init(cfg)
params on device; CUDA kernels for matmul/softmax/GELU/LN
adam_init (with linear_lr.decay)

**Optional param broadcast**
if world_size > 1:
rank 0 calls mpi_broadcast_parameters
all ranks start from identical weights

**Training loop**
**1. Local batch**
dataloader.next_batch on each rank
**2. Forward (CUDA)**
gpt_forward_with_loss: device logits, scalar loss
**3. Backward**
backward(loss) on GPU, clear activations
**4. Optional grad allreduce**
if world_size > 1: mpi_allreduce_grads
sums gradients across ranks
**5. Optimizer**
adam_step (+ optional grad clipping)
adam_zero_grad
master rank prints loss

**Checkpoint**
master rank: save_weights("trained_weights.bin")

**Text generation (master rank)**
sliding window tokens on device
gpt_forward_logits (CUDA)
greedy_select_next_token
append token; shift window
stop on EOS or length cap
decode IDs $\to$ text and print

## Stage 3

**MPI + CUDA init**
MPI_Init
RNG: tensor.set_seed(seed + rank)
CUDA devices: cudaGetDeviceCount,
per-rank device: cudaSetDevice(rank % n.devices)

**Tokenizer on rank 0**
Tokenizer("../data/dummy.txt")
tokenizer_extract $\to$ tokenizer_encode

**Tokenizer broadcast**
broadcast success flag; tokenizer_broadcast
all ranks share vocab & encoded data
tokenizer_pad_to(global_batch_size)

**Rank-aware DataLoader**
per-rank batch: $B_{rank} = B_{global}/world\_size$
DataLoader.init_with_tokenizer(tokenizer, rank, world_size)
sharded sampling along token stream

**Model init on GPU**
gpt_init(cfg) (weights on device)
gpt_set_distributed for MPI data parallel
CUDA kernels for matmul/softmax/GELU/LN/xent
adam_init (with linear_lr.decay)

**Initial parameter sync**
if world_size > 1:
rank 0 calls mpi_broadcast_parameters
all ranks start from identical GPU weights

**Training loop**
**1. Sharded batch**
dataloader.next_batch uses rank, world_size
**2. Forward on GPU**
gpt_forward_with_loss: logits, scalar loss
**3. Backward**
backward(loss), clear activations
**4. Gradient allreduce**
mpi_allreduce_grads sums grads across ranks
**5. Optimizer**
adam_step (optional grad clipping)
adam_zero_grad
**6. Loss reduction & logging**
mpi_allreduce_loss (average loss)
master rank prints epoch loss

**Checkpoint**
master rank: save_weights("trained_weights.bin")

**Text generation (master rank)**
sliding window token IDs on GPU
gpt_forward_logits
greedy decode: greedy_select_next_token
append token, update window
stop on EOS or length limit
decode IDs $\to$ text and print

# Hyperparameters

**GPT model hyperparameters:**

- Block size: 256
- Layers: 6
- Heads: 6
- Embedding size: 384
- Dropout: *N/A*
- Learning rate: 3e-4f

**Data-loading stats:**

- Batch size: 16
- Sequence length: 256
- Epochs: 50

**Toy dataset stats:**

- 290,599 bytes
- Extracted 52,600 tokens
- Vocab size: 172

*Note: the dataset is small and highly repetitive

# Training Results

- We ran host MPI with 4 processes
  - Why? The Intel(R) Core(TM) i7-7820HQ CPU only has 4 cores (2 threads per core)
- All versions used the same hyperparameters

| Versions | Training Time |
|---|---|
| Serial | 15264s (254.4 minutes) |
| MPI (Two nodes) | 1287s (21.5 minutes) |
| CUDA | 429s (7.2 minutes) |
| MPI+CUDA (Two nodes) | 226s (3.8 minutes) |

# Inference Results

| Versions | Inference Time |
|---|---|
| Serial | 113s (1.9 minutes) |
| MPI (Two nodes) | 112s (1.9 minutes) |
| CUDA | 8.0s (0.13 minutes) |
| MPI+CUDA (Two nodes) | 7.8s (0.13 minutes) |

# Conclusions

- We found significant speedup,

    - The MPI implementation yielded an immense speedup when utilizing both the CPU and CUDA (GPU)

        resources.

- You should,

    - Keep CUDA kernels simple

    - Minimize unnecessary data transfers between the host and the device

- Avoid,

    - Blocking calls for synchronization - can cause deadlocks or unnecessary waiting

    - Sending many small chunks of data rather than large chunks for communication calls (MPI)

# Future Direction

- Add NCCL (NVIDIA Collective Communications Library)
  - The standard for distributed training and inference across multiple GPUs.
  - Combines MPI and CUDA routines into a single framework.
  - Used to train actual AI models.
- Can also increase:
  - Hyperparameters
  - # of processes and/or machines
  - Data size