# Project Documentation: Accelerating PyTorch with Parallelization

Noah Snodgrass, Minh (Ben) Pham, Jerry Li, Gokul Nookula, Samuel Uribe, Jon Darius

## Overview

This project focuses on understanding PyTorch, a powerful library for machine learning, and exploring ways to accelerate its computations through parallelization. Our goal is to identify performance bottlenecks in PyTorch's operations and implement efficient parallel processing techniques to optimize its functionality.

We base our work on CTorch, an open-source project that re-implements portions of PyTorch in C++, with a focus on GPU acceleration. By leveraging the concepts and architecture of CTorch, we aim to further enhance the speed and efficiency of PyTorch's matrix multiplication, backpropagation, tensors, linear functions, etc.

## Motivation

Deep learning frameworks like PyTorch rely heavily on computational performance for training and inference tasks. While PyTorch already supports GPU acceleration and optimized backends, there are opportunities for further optimizations:

1. Efficient Memory Utilization: Minimize memory overhead and improve cache efficiency
2. Parallel Execution: Leverage advanced parallel programming techniques to maximize throughput
3. Custom Optimizations: Implement task-specific enhancements to critical operations like matrix multiplications and convolutions

This project provides an excellent opportunity to delve into the internals of PyTorch, learn advanced parallel programming techniques, as well as basic machine learning.

## Key Objectives

- Understand PyTorch Fundamentals
- Optimize Computation:
    - Identify bottlenecks in CTorch or PyTorch operations
    - Implement parallel processing techniques to accelerate operations (stated before)

● Applying parallel programming to machine learning

# Methods of Parallelization Used

- Multithreading for Tensor Operations:
- Context: Operations on tensors, such as addition, subtraction, and scalar manipulations, are parallelized to utilize multi-core CPUs.

- Implementation: Multithreading is applied using the std::thread library. For example:

- Addition with Scalar:
  - Threads are divided into chunks, and each thread processes a segment of the tensor's data.
    - Example: In tensor_example += 1, each individual thread handles a portion of the tensor's elements

- Addition/Subtraction of Tensors:
  - Threads run independently and join at the end, similar to global sum problem

- The mnist.cpp file demonstrates using std::async to normalize a dataset of images.
  - The dataset is divided among 4 threads (numThreads = 4), where each thread normalizes a subset of images
  - Example:
    - Each thread processes a subsetSize of images independently.
    - Synchronization is achieved using future.get() to wait for all threads to complete their tasks

- Performance:
  - Timing is recorded using chrono to measure the efficiency of normalization.
  - This approach enables faster processing compared to a single-threaded implementation.

# Efficient Data Access with Array and Tensor Classes:

- Context: The Array and Tensor classes form the foundation for operations on multidimensional data
- Features:
  - Strides Calculation: Enables efficient memory access patterns for operations like flattening and transposing
  - Thread-Safe Access: The Array class provides thread-safe access mechanisms, ensuring no race conditions occur during multithreaded computations
  - The matrix multiplication operator (operator*) supports 2D tensors

- Integration:

- ○ The Tensor class integrates with Array to perform element-wise and matrix operations using parallelized methods, used in the Tensor class as well
- ○ The matrix multiplication operator (operator*) supports 2D tensors and ensures compatibility before computation.

# Design Considerations
- Scalability:
  - ○ The number of threads can be adjusted dynamically based on the hardware capabilities and workload size
  - ○ The modular design allows adding more parallelized operations with minimal refactoring
- Error Handling:
  - ○ The code includes shape validation and dimension checks to ensure operations are well-defined

- The mnist.cpp file includes timing mechanisms to record the duration of normalization tasks which provides insights into the speedup achieved from parallelization

# Conclusion
- The parallelization implemented in the provided code significantly enhances computational efficiency, particularly for large datasets and intensive operations. By leveraging multithreading and asynchronous programming, the code ensures optimal utilization of system resources while maintaining correctness and scalability. These optimizations are a step forward in creating high-performance frameworks for machine learning and data processing.