



Name: Salma Ahmed Sherif

Date: March 28th, 2025

As someone deeply passionate and well-experienced about low-level programming and modern AI, I am excited to work at the intersection of C++ and large language models (LLMs). Please find a report that details my work through the evaluation task.

Table of Contents

1. Introduction

- Initial Transformer Implementation
- Challenges with Clad Integration
- Simplification to a Linear Model
- Testing and Validation

2. Methodology Overview

2.1. Phase 1: Transformer Implementation with Eigen

- 2.1.1. Token Embedding
- 2.1.2. Scaled Dot-Product Attention
- 2.1.3. Transformer Block
- 2.1.4. Training Loop

2.2. Phase 2: Simplification for Clad Integration

- 2.2.1. Challenges with Clad
- 2.2.2. Simplified Linear Model (Code Structure & Clad Integration)
- 2.2.3. Data Structures Used

2.3. Unit Tests and Validation

- Forward Pass Validation
- Gradient Verification
- Handling Discontinuities
- Jacobian Validation

3. Literature Review

- 3.1. Foundational Architecture (Vaswani et al., 2017)
- 3.2. Low-Level Implementations (llm.cpp)

- 3.3. Automatic Differentiation (Clad)

4. High-Level Implementation Plan

- 4.1. Transformer Architecture with Eigen (`train_transformer.cpp`)
 - 4.1.1. Creative Innovations (Eigen-Centric Design, Dynamic Batching)
- 4.2. Simplified Linear Model with Clad (`train_multiple_neurons.cpp`)
 - 4.2.1. Creative Innovations (Array-Centric Simplification, Jacobian-Centric Training)

5. Commands and Instructions

- Compilation and Execution for Transformer
- Compilation and Execution for Linear Model

6. Pull Requests and Issues

- Submitted Issues (#1306, #1313, #1320)
- Submitted Pull Requests (#1297)

References

1. Introduction

This report outlines the development of a minimal Large Language Model (LLM) in C++, starting with a transformer architecture and transitioning to a simplified linear model to integrate Clad for efficient automatic differentiation. The workflow proceeded as follows:

1. Initial Transformer Implementation:

- A transformer block with multi-head attention and feed-forward layers was implemented using Eigen's matrix library (`train_transformer.cpp`).

2. Challenges with Clad Integration:

- Eigen's matrix operations (`MatrixXd` , `VectorXd`) conflicted with Clad's differentiation requirements, which favor primitive arrays and fixed-size matrices.
- **Critical Code Sections:**
 - **Manual Gradients:** In `train_transformer.cpp` , gradients for weights like WQ,WK,WVWQ,WK,WV were computed using Eigen's matrix algebra.

```
// Manual gradient calculation (TransformerBlock::train_step)
VectorXd d_pred = pred - target;
MatrixXd dW2 = d_pred * hidden.transpose();
```

3. Simplification to a Linear Model:

- A minimal multi-output linear model (`training_multiple_neurons.cpp`) was developed to align with Clad's function signatures and data structures.
- **Rationale:** Clad's `clad::matrix` type and Jacobian computation required reworking the transformer's matrix-heavy architecture into a parameter array-based design.

4. Testing and Validation:

- The simplified model achieved stable training loss reduction over 100 epochs (e.g., final loss: `0.0113776`).
 - Clad successfully computed gradients for all independent parameters , replacing manual backpropagation.
-

2. Methodology Overview

This section details the evolution of the implementation, starting with a transformer model using Eigen for matrix operations and transitioning to a simplified linear model integrated with Clad for automatic differentiation. The methodology is divided into **two** phases:

2.1. Phase 1: Transformer Implementation with Eigen

The core of the methodology is based on implementing a **transformer model** in C++ to build a minimal Large Language Model (LLM). The key components are token embedding, multi-head self-attention, feed-forward layers, and a simple training loop using gradient descent.

2.1.1. Token Embedding:

In the code, the input tokens are **embedded** into fixed-size vectors of `EMBEDDING_SIZE = 16` . This is crucial as it transforms discrete tokens into continuous vectors that the model can understand. For simplicity, the embeddings are generated randomly using the Eigen matrix library. In practice, embedding vectors would be learned during training.

Code Excerpt:

```
// Random embedding generation for input tokens
std::vector<MatrixXd> inputs = { MatrixXd::Random(EMBEDDING_SIZE, 1),
MatrixXd::Random(EMBEDDING_SIZE, 1) };
```

2.1.2 Scaled Dot-Product Attention:

The core of the transformer is the **scaled dot-product attention** mechanism. This mechanism compares each token (query `Q`) to every other token (key `K`) in the sequence and computes how much focus each token should give to the others. The result is a set of attention scores that influence the final output.

Code Excerpt:

```
// Scaled Dot-Product Attention
MatrixXd scaled_dot_product_attention(const MatrixXd& Q, const MatrixXd& K, const
```

```

MatrixXd& V) {
    MatrixXd scores = (Q * K.transpose()) / sqrt(K.cols()); // Scale the dot
product
    MatrixXd attention = scores.array().exp(); // Apply the softmax
    return attention * V; // Weighted sum of the values (V)
}

```

The **Eigen Matrix** is chosen here because it provides an efficient way to handle matrix operations and is highly optimized for performance. Using Eigen makes the matrix multiplications for dot products and weight updates straightforward and fast.

2.1.3. Transformer Block:

A transformer block is implemented that consists of the attention mechanism and a feed-forward network. It includes the weight matrices `Wq`, `Wk`, `Wv`, and `Wo`, which are used for queries, keys, values, and the output of the attention mechanism, respectively. The block also contains hidden layers, `W1` and `W2`, which are fully connected layers for processing the attention output.

Code Excerpt:

```

// Transformer Block
class TransformerBlock {
public:
    MatrixXd Wq, Wk, Wv, Wo;
    MatrixXd W1, W2;
    VectorXd b1, b2;

    TransformerBlock(int dim, int heads) {
        // Initialize weights and biases randomly
        Wq = MatrixXd::Random(dim, dim);
        Wk = MatrixXd::Random(dim, dim);
        Wv = MatrixXd::Random(dim, dim);
        Wo = MatrixXd::Random(dim, dim);
        W1 = MatrixXd::Random(HIDDEN_SIZE, dim);
        W2 = MatrixXd::Random(dim, HIDDEN_SIZE);
        b1 = VectorXd::Random(HIDDEN_SIZE);
        b2 = VectorXd::Random(dim);
    }

    // Forward pass
    MatrixXd forward(const MatrixXd& input) {
        MatrixXd Q = Wq * input;
        MatrixXd K = Wk * input;
        MatrixXd V = Wv * input;
        MatrixXd attention_output = scaled_dot_product_attention(Q, K, V);
        MatrixXd output = Wo * attention_output;

        VectorXd hidden = (W1 * output.col(0) + b1).array().tanh();
    }
}

```

```

        return W2 * hidden + b2;
    }
};

```

The **MatrixXd** type is used here as it represents a matrix with double precision, suitable for large matrix operations, and **VectorXd** is used for handling vectors in the feed-forward layers. These data structures are well-suited to C++ and Eigen's efficient handling of mathematical operations.

2.1.4. Training Loop:

The **training loop** runs for multiple epochs, where each epoch consists of calculating the loss, performing backpropagation, and updating the model weights using gradient descent. In this implementation, the loss is computed using **Mean Squared Error (MSE)**, and the weights are updated manually based on the gradients.

Code Excerpt:

```

// Training loop
void train(TransformerBlock& model, const std::vector<MatrixXd>& data, const
std::vector<VectorXd>& labels) {
    for (int epoch = 0; epoch < EPOCHS; ++epoch) {
        double total_loss = 0.0;
        for (size_t i = 0; i < data.size(); i++) {
            model.train_step(data[i], labels[i]); // Perform a training step
        }
    }
}

```

In this function, the `train_step` method is invoked for each data instance, which computes the forward pass, calculates the loss, and updates the weights.

2.2. Phase 2: Simplification for Clad Integration

2.2.1 CHALLENGES WITH CLAD:

Clad's AD system operates on primitive arrays and fixed-size matrices, conflicting with Eigen's dynamic `MatrixXd`. Function signatures for `'clad::jacobian'` required explicit parameter ordering, incompatible with Eigen's object-oriented design.

- **Simplified Linear Model:**

A minimal multi-output linear model (`training_multiple_neurons.cpp`) was developed to align with Clad's requirements:

The provided code demonstrates a multi-output linear model that computes predictions using two neurons, each with its own weights and biases. The core structure consists of the following key components:

- *Forward Function:* This computes the predictions (y_1 and y_2) using a simple linear transformation of the inputs. Each neuron computes a weighted sum of the inputs followed by the addition of a bias term.
- *Loss Function:* The Mean Squared Error (MSE) is used to compute the loss between the predicted outputs and the target outputs.
- *Gradient Calculation:* This was initially implemented manually in the code, where gradients were computed with respect to each weight and bias.
- *Integration of Clad for Differentiation:* The major enhancement is the integration of Clad for automatic differentiation, specifically calculating the Jacobian matrix, which allows the efficient calculation of gradients.

2.2.2. KEY CODE EXCERPTS AND THEIR DESCRIPTION

Forward Function (Linear Model)

The forward function calculates the predictions for two output neurons based on input features and weights. The weights W_0, W_1, W_2, W_3 are used to linearly transform the inputs x_0, x_1 , while the biases b_1, b_2 are added to the results for each neuron.

```
void forward(  
    double W0, double W1, double b1,  
    double W2, double W3, double b2,  
    double x0, double x1, double* y) // Use an array for multiple outputs  
{  
    y[0] = W0 * x0 + W1 * x1 + b1; // First neuron  
    y[1] = W2 * x0 + W3 * x1 + b2; // Second neuron  
}
```

This function forms the basis of the model, applying linear transformations and storing the results in the y array.

Loss Function (Mean Squared Error)

The loss function is the Mean Squared Error (MSE), which calculates the difference between predicted outputs and actual targets.

```
double loss(const double* y, const double* target) {  
    double loss_val = 0;  
    for (int i = 0; i < OUTPUT_SIZE; i++) {  
        loss_val += 0.5 * (y[i] - target[i]) * (y[i] - target[i]);  
    }  
}
```

```

    }
    return loss_val;
}

```

MSE is a simple yet effective loss function for regression tasks, which this linear model appears to perform.

Clad Integration for Automatic Differentiation

The integration of Clad enables the automatic computation of gradients using its jacobian function. The Jacobian matrix represents the partial derivatives of the model's outputs with respect to its parameters (weights and biases).

```

// Compute Jacobian using Clad
clad::matrix<double> d_output(2, 10); // Initialize Jacobian with correct
dimensions (2 outputs × 10 parameters)
double temp_y[OUTPUT_SIZE] = {0.0, 0.0};

// Execute jacobian with correct parameter order
clad::jacobian(forward).execute(
    W[0], W[1], b[0], W[2], W[3], b[1],
    x0, x1, temp_y, &d_output
);

```

Clad's jacobian function computes the partial derivatives for each parameter, allowing gradient-based optimization. This change significantly optimizes the differentiation process and makes the model more extensible.

Gradient Computation and Weight Update

The gradients for each parameter are computed using the Jacobian matrix, and the weights are updated using gradient descent. For each epoch, the gradients are accumulated for each parameter and averaged across the batch.

```

// Apply gradient descent update
for (int k = 0; k < 4; k++) W[k] -= LEARNING_RATE * total_dW[k];
for (int k = 0; k < 2; k++) b[k] -= LEARNING_RATE * total_db[k];

```

This ensures that the model iterates towards minimizing the loss by adjusting the weights and biases based on the computed gradients.

2.2.3. DATA STRUCTURES USED

- Clad's Matrix Class (clad::matrix): This is used for storing the Jacobian matrix, which holds the derivatives of the model's outputs with respect to its parameters.

- **Standard Arrays for Weights and Biases:** The weights $W[]$ and biases $b[]$ are stored as simple arrays because the number of parameters is small and they can be efficiently updated during training.
- **Vectors for Input and Target Data:** The training data is represented using standard C++ vectors (`std::vector`) for flexibility.

2.3. Unit Tests and Validation

To ensure the correctness and robustness of the implementation, rigorous testing was performed at each stage of development. The following validation strategies were employed:

1. **Forward Pass Validation:**
 - **Objective:** Verify that the model's forward pass produces mathematically correct outputs.
 2. **Gradient Verification:**
 - **Objective:** Ensure Clad-computed gradients match ground-truth values.
 3. **Handling of Discontinuities:**
 - **Objective:** Test edge cases in activation functions (e.g., ReLU at $x=0$).
 4. **Jacobian Validation:**
 - **Objective:** Confirm Clad's Jacobian matrices align with theoretical expectations.
-

3. Literature Review

This section highlights foundational works and tools critical to developing efficient, low-level LLMs in C++, with a focus on transformer architectures, performance-optimized implementations, and automatic differentiation.

3.1. Foundational Architecture

- **"Attention is All You Need" (Vaswani et al., 2017):**
Introduced the transformer architecture, which underpins modern LLMs. Its self-attention mechanism and encoder-decoder structure directly informed the design of the multi-head attention and feed-forward layers in this project.

3.2. Low-Level Implementations

- **llm.cpp (Zhangpiu, 2023):**
Provided a minimalist, open-source implementation of transformer training in C++. Its use of raw pointers and batched gradient updates influenced the project's initial design for lightweight training loops.

3.3. Automatic Differentiation

- **Clad (Vassilev et al., 2021):**
A Clang-based automatic differentiation (AD) tool for C++. Its Jacobian computation capabilities replaced manual gradient calculations in the simplified linear model, ensuring maintainability and reducing error-prone code. Clad’s compatibility with procedural C++ made it preferable to template-based alternatives like ADOL-C.
Summary:

Resource	Role in Implementation
Vaswani et al. (2017)	Theoretical foundation for attention and transformer blocks.
llm.cpp	Inspired minimalist training loop and matrix handling.
Clad	Enabled automated gradient computation for linear models.

4. High-Level Implementation Plan

This section outlines the architectural design and innovative strategies employed in two distinct implementations: a transformer model using Eigen for matrix operations (`train_transformer.cpp`) and a simplified linear model integrated with Clad for automatic differentiation (`train_multiple_neurons.cpp`). The evolution from a complex transformer to a minimal linear model reflects a pragmatic balance between theoretical rigor and practical differentiation needs.

4.1. Transformer Architecture with Eigen (`train_transformer.cpp`)

The initial implementation focused on a minimalist transformer model, prioritizing efficiency and scalability through low-level control. Key components include:

- **Token Embedding:** Input tokens are mapped to dense vectors using dynamically sized matrices from the Eigen library, enabling flexible dimensionality while maintaining computational efficiency.
- **Scaled Dot-Product Attention:** The core attention mechanism computes pairwise token dependencies using optimized matrix multiplications and scaling, leveraging Eigen’s SIMD-enabled operations for performance.
- **Feed-Forward Network:** A two-layer neural network processes attention outputs, incorporating non-linear activations (e.g., hyperbolic tangent) to enhance representational capacity.
- **Manual Gradient Computation:** Gradients for weight matrices (e.g., query, key, value projections) were derived manually using matrix calculus, exploiting Eigen’s algebraic expressiveness for clarity.
- **Random Initialization:** Weights are initialized randomly, which is standard for neural networks to break symmetry and allow for effective learning.

4.1.1.CREATIVE INNOVATIONS:

- **Eigen-Centric Design:** By relying on Eigen's `MatrixXd` and `VectorXd`, the implementation achieved high-performance matrix operations without external dependencies, ensuring portability and efficiency.
- **Dynamic Batching:** The training loop processes inputs as batches of Eigen matrices, enabling parallelism while retaining memory locality.

4.2. Simplified Linear Model with Clad (`train_multiple_neurons.cpp`)

To integrate Clad's automatic differentiation, the transformer architecture was distilled into a multi-output linear regression model, prioritizing compatibility with Clad's procedural workflow:

- **Linear Forward Pass:** The model computes predictions via linear transformations of inputs using primitive arrays, eliminating Eigen's object-oriented overhead and aligning with Clad's function signature requirements.
- **Clad-Driven Differentiation:** Clad's `jacobian` function automates gradient computation by tracing parameter dependencies in the forward pass, replacing error-prone manual derivatives with a maintainable, compiler-aided approach.
- **Batch Gradient Descent:** Gradients are averaged across mini-batches to stabilize training, with updates applied directly to array-based parameters for simplicity.

4.2.1. CREATIVE INNOVATIONS:

- **Array-Centric Simplification:** Replacing Eigen matrices with primitive arrays reduced complexity, enabling seamless integration with Clad's differentiation engine while retaining computational efficiency.
- **Jacobian-Centric Training:** By structuring the forward pass to expose all parameters as scalar arguments, Clad's Jacobian computation could be directly applied, a novel adaptation for linear models in C++.

Comparative Design Choices

Aspect	Transformer (Eigen)	Simplified Model (Clad)
Matrix Operations	Eigen's dynamic <code>MatrixXd</code> for flexibility	Primitive arrays for Clad compatibility
Differentiation	Manual gradient computation	Automated via Clad's Jacobian
Data Structures	Object-oriented (Eigen classes)	Procedural (arrays, Clad matrices)
Key Advantage	High-performance attention mechanisms	Streamlined differentiation and maintainability

5. Commands and Instructions

The code was compiled and executed using the following commands:

1. Transformer Implementation (`train_transformer.cpp`)

```
sudo apt-get install libeigen3-dev

clang++ -O3 -std=c++20 -fopenmp train_transformer.cpp -o train_transformer

./train_transformer
```

2. Linear Model with Clad (`train_multiple_neurons.cpp`)

```
clang++ -std=c++17 -O0 -g \
-I/root/llvm-project/llvm/tools/clad/include \
-I/usr/include/eigen3 \
-fplugin=/root/llvm-project/llvm/tools/clad/build/lib/clad.so \
train_multiple_neurons.cpp -o multiple

./multiple
```

6. Pull Requests and Issues

To improve Clad's robustness and address edge cases encountered during development, the following contributions were made:

Submitted Issues

1. [Issue #1306](#):

- **Description:** I found the problem to be that users get a cryptic C++ compiler error later, rather than a clear Clad diagnostic at that scenario.

2. [Issue #1313](#):

- **Description:** I found the problem to be a segmentation fault instead of a clear Clad diagnostic or handling the case in the code.

3. [Issue #1320](#):

- **Description:** In the context of training loops, I was experimenting with exponential functions

when I found a 'core dump'

Submitted Pull Requests

1. PR #1297:

In addressing issue #1265 in the Clad project, where forward-declared pullback functions, which are not utilized, could lead to compilation errors due to the absence of function definitions. To resolve this, I added empty function bodies for these unused pullback functions. Building clad successfully after the edits and creating a corresponding test file were part of the process.

References

1. Vaswani, A., et al. (2017). "Attention is All You Need." *NeurIPS*.
2. Zhangpiu. (2023). llm.cpp. *GitHub Repository*. <https://github.com/zhangpiu/llm.cpp>
3. Vassilev, V., et al. (2021). "Clad: Automatic Differentiation for C++." *CERN*.