

High Performance Computer Architectures Practical Course - Exercise 3 -

Tutorium 1

David Jordan (6260776)
Florian Rüffer (7454628)
Michael Sanjatin (7485765)

May 9, 2023

1 Neural Networks and SIMD

1

The task is to adapt the feedForward function in the MLP.cpp file. In order for this to work the affineTransform, SoftMax and MatMul2D1D had to be adapted to the self-defined fvec datatype.

```
1  std::vector<fvec> MLPMath::MatMul2D1D(std::vector<std
    ::vector<float>>& A, std::vector<fvec>& B) {
2
3      std::size_t nArows = A.size();
4      std::size_t nAcolumns = A[0].size();
5      std::size_t nBrows = B.size();
6
7      // initialize matrix C to all zero values
8      std::vector<fvec> C(nArows);
9
10     // check if A and B are commensurate and multiply
11     if (nAcolumns == nBrows) {
12         // -----
13         // loop over all rows of the matrix A
14         for(std::size_t j = 0; j < nArows; j++){
15             // declare an fvec to store the
16             // intermediate results in
17             fvec tempo = fvec(0.0f);
18             // store the j-th row of A in another
19             // temporary variable
20             std::vector<float>& tmp = A[j];
21             // looping over the entries and
22             // multiplying
23             for (std::size_t l = 0; l < nAcolumns; l
24                 ++){
25                 // utilizing the fvec multiplication
26                 // for SIMD
27                 tempo += tmp[l] * B[l];
28             }
29             C[j] = tempo;
30         }
31         // -----
32         return C;
33     }
34     else {
35         std::cerr << "Matrix indexes incompatible for
36             multiplication. A columns =" << nAcolumns
```

```

32         << "B rows =" << nBrows << std::endl;
33         exit(EXIT_FAILURE);
34     }
35 }
36 // The function applies the softmax activation
37 // function to an input vector of type fvec
38 std::vector<fvec> MLPMath::applySoftmax(std::vector<
39 fvec>& input) {
40     int inputSize = (int)input.size();
41     // The output vector is initialised to be of the
42     // same size as the input vector
43     std::vector<fvec> output(inputSize);
44     // A vector to store the exponential values for
45     // each input vector
46     std::vector<float> exponentialValue(inputSize);
47     // Stores the sum of the exp values
48     float exponentialSum = 0.0f;
49     // -----
50     // Each column of the input matrix is looped over
51     for(int i = 0; i < fvecLen; i++){
52         // determines the current max value in the
53         // selected column
54         float maximumV = input[0][i];
55         for (std::size_t j = 1; j< inputSize; j++){
56             if (maximumV < input[j][i]){
57                 maximumV = input[j][i];
58             }
59         }
60         // Determine the exponential values, first
61         // clearing the values
62         exponentialSum = 0.0f;
63         exponentialValue.clear();
64         exponentialValue.resize(inputSize);
65         // For each input vector in the selected
66         // column calculate the exponential value
67         for (std::size_t k = 0; k<inputSize; k++){
68             float temp = exp(input[k][i] - maximumV);
69             exponentialValue[k] = temp;
70             exponentialSum += temp;
71         }
72     }

```

```

70         // Each exp value is divided by the sum of all
           exponential values to calculate the
           softmax output
71     for (std::size_t l = 0; l<inputSize; l++){
72         output[l][i] = exponentialValue[l] /
           exponentialSum;
73     }
74 }
75 // -----
76
77 return output;
78 }
79
80 std::vector<fvec> MLPMath::affineTransform(std::vector
<std::vector<float>>& weight, std::vector<fvec>&
input, std::vector<float>& bias) {
81
82     std::vector<fvec> result(bias.size());
83     std::vector<fvec> term1 = MLPMath::MatMul2D1D(
        weight, input);
84
85     for (std::size_t i = 0; i < bias.size(); i++) {
86         result[i] = term1[i] + bias[i];
87     }
88
89     return result;
90 }

```

File 1: SIMDed math functions

The functions now utilize the SIMDed arithmetic operators of the fvec class.

2

The task for this exercise is to add SIMD (Single Instruction, Multiple Data) support to our neural network functions. We will start of with the simd-ized version of the ReLU activation function.

The first noticable detail is the changed input. We input type is "fvec". Fvec is a SIMD-packed set of values, to be exact 4 packed 32-bit floating point numbers. Once again we create a result vector, the same size as the input. Subsequently we loop for 1.) every element of fvec (here this is 4) and 2.) the number of elements of the input. As a result, we loop over every element of the input.

Finally, we check the conditions associated with the ReLU activation function for every of those elements. The defintion of ReLU is provided in the Appendix

(5.1)

```
1      std::vector<fvec> MLPMath::applyReLU(std::vector<
      fvec>& input) {
2
3      std::vector<fvec> result(input.size());
4
5      for (std::size_t iv = 0; iv < fvecLen; iv++) {
6          for (std::size_t i = 0; i < input.size(); i++)
              {
7              if (input[i][iv] > 0.0f) {
8                  result[i][iv] = input[i][iv];
9              }
10             else {
11                 result[i][iv] = 0.0f;
12             }
13         }
14     }
15
16     return result;
17 }
```

File 2: SIMD-ized ReLU

Another essential part is the backPropActivation function. Here we provide just the code snippet, which specifically implements ReLU.

In this case there are not many significant changes. We want to ensure, that this function will execute properly if provided with a fvec vector input. For this to work, line 4 is the most essential one to understand. Here the fvec-type overloads the ">"-operator. To understand the underlying functionality, please view the following example.

"rawOutput" contains the elements: [1.0, -1.0, 2.0, -0.5]

"zero" contains the elements [0.0, 0.0, 0.0, 0.0]

Overloading will produce a result like this: [1.0 > 0.0, -1.0 > 0.0, 2.0 > 0.0, -0.5 > 0]

This checks the condition for every element and can be abbreviated into the following form: [0xFFFFFFFF, 0x00000000, 0xFFFFFFFF, 0x00000000] (1 - True & 0 - False)

```
1      case 1:{
2          fvec zero = fvec(0.0f);
3          for (std::size_t i = 0; i < rDelta.size(); i
4              ++){
5              rDelta[i] = (rawOutput[i] > zero) * rDelta
6                  [i];
7          }
8          break;
9      }
```

File 3: backPropActivation

Finally we need to incorporate our SIMD-ized ReLU function into the apply-Activation function (Line 10-14).

```
1      std::vector<fvec> MLPNet::applyActivation(std::
2          vector<fvec>& input) {
3
4      std::vector<fvec> output(input.size());
5      switch (activationType_) {
6          case 0: {// TanH
7              output = MLPMath::applyTanH(input);
8              return output;
9              break;
10         }
11         case 1: {// SIMD-ized ReLU
12             output = MLPMath::applyReLU(input);
13             return output;
14             break;
15         }
16         default: {
17             return output;
18             break;
19         }
20     }
```

File 4: applyActivation

3

When the program is run on the "mnisttrain.csv" and "mnisttest.csv" data, and the required dimensions of 10000 training / test examples, a batch size of 40 with 10 epochs and the RELU activation function the following output is recorded:

```
[100%] Built target HPC_MLP
hpc@hpc-vm:~/Downloads/simd/simd_ex/HPC-MLP-SIMD-exer/bu
MLP topology : [
( 784 , 800 ),
( 800 , 10 ) ].
Number of Training examples : 10000
Batch Size : 40
Number of Epochs : 10

Correct = 5824 , Wrong = 4176
Correct = 6933 , Wrong = 3067
Correct = 7303 , Wrong = 2697
Correct = 7361 , Wrong = 2639
Correct = 6900 , Wrong = 3100
Correct = 7316 , Wrong = 2684
Correct = 7154 , Wrong = 2846
Correct = 7205 , Wrong = 2795
Correct = 7411 , Wrong = 2589
Correct = 7353 , Wrong = 2647
Training Accuracy in Epoch 1 = 58.24
Training Accuracy in Epoch 2 = 69.33
Training Accuracy in Epoch 3 = 73.03
Training Accuracy in Epoch 4 = 73.61
Training Accuracy in Epoch 5 = 69
Training Accuracy in Epoch 6 = 73.16
Training Accuracy in Epoch 7 = 71.54
Training Accuracy in Epoch 8 = 72.05
Training Accuracy in Epoch 9 = 74.11
Training Accuracy in Epoch 10 = 73.53
Number of Testing examples ( 0 to skip testing) : 10000
Correct = 7180 , Wrong = 2820
```

Figure 1: Output

The peak training accuracy of 74.11 in epoch 9 and a testing accuracy of 71.80 is recorded when we run the program on the data.

2 Matrix

In this task we need to speed up Matrix calculations using SIMD. As displayed below, we want to go from matrix a to matrix c .

$$a = \begin{bmatrix} a_0 & a_1 & a_2 \\ a_3 & a_4 & a_5 \\ a_6 & a_7 & a_8 \end{bmatrix}$$
$$c = \begin{bmatrix} \sqrt{a_0} & \sqrt{a_1} & \sqrt{a_2} \\ \sqrt{a_3} & \sqrt{a_4} & \sqrt{a_5} \\ \sqrt{a_6} & \sqrt{a_7} & \sqrt{a_8} \end{bmatrix}$$

First of all we need to change the input & output matrices (output for scalar computation stays the same). We do this by changing the memory alignment of the values of the matrices. This is a requirement for SIMD-ized operations.

```

1 float a[N][N] __attribute__((aligned(16))); //
  input array
2 float c[N][N]; // output array for scalar computations
3 float c_simd[N][N] __attribute__((aligned(16))); //
  output array for SIMD computations

```

File 5: Matrix.cpp

The loops for the computation part remain mostly unchanged to the scalar version. The inner loop runs N -times but our loop variable increases by 4, because each vector has 4 elements and those can be calculated in parallel. We loop for $N_{iter} - 1$ times, to neglect memory reading time. Subsequently we loop over matrix a to transform it into matrix c ($N \times N$ for rows and columns). The distinct section ranges from line 5 to line 7. We define two new vectors, called $aVec$ and $cVec$, which suppose to be the vector counter parts to the scalar versions of matrix a and matrix c . At the beginning those were initialized with float values, so we use the `reinterpret_cast` command to allow treating a variables memory representation as another type. In this case we substitute float for `fvec`. Finally we just plug those values into the template function to calculate the root (Line 7).

```

1 TStopwatch timerSIMD;
2 for( int ii = 0; ii < NIter; ii++ )
3     for( int i = 0; i < N; i++ ) {
4         for( int j = 0; j < N; j++ ) {
5             fvec &aVec = reinterpret_cast<fvec*>(a[i][
              j]);
6             fvec &cVec = reinterpret_cast<fvec*>(
              c_simd[i][j]);
7             cVec = f(aVec);
8         }
9     }
10 timerSIMD.Stop();

```

File 6: Matrix.cpp

We can compute 4 values in parallel and therefore can expect a speed-up factor of 4. Due to running environment deviation, the speed-up factor varies.

```

Time scalar: 355.03 ms
Time SIMD: 405.338 ms, speed up 0.875886
SIMD and scalar results are the same.

```

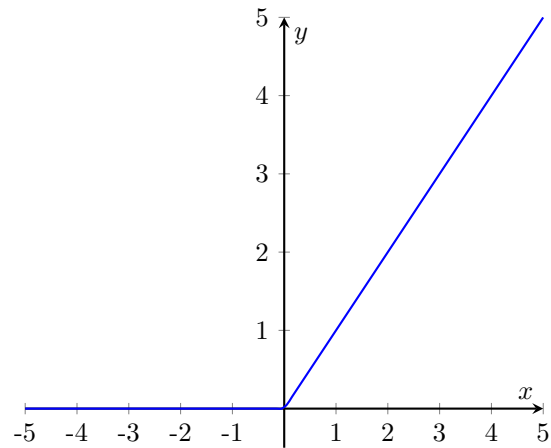
Figure 2: Output

3 Quadratic Equation

4 CheckSum

5 Appendix

5.1 ReLU



$$\text{ReLU}(x) = \begin{cases} x, & \text{if } x \geq 0 \\ 0, & \text{otherwise} \end{cases}$$



Figure 3: Add caption