

High Performance Computer Architectures

Practical Course

- Exercise 3 -

May 3, 2023

Neural Networks and SIMD

This week, there are several tasks on SIMD instructions in C++. The MLP implementation of the previous week should be SIMDized over the batches. Furthermore, there are old exercises to solve.

We will now use SIMD in our neural network package. You might have noticed that all the examples in a batch can be run through the neural network independently. It is only at the end of a batch that we sum up all the gradients and update the weights. So we can pack together `fvecLen` number of training examples of a batch into a SIMD vector and run them through the neural network together.

Problem

Follow the comments in the code to complete the SIMD-ized neural network:

1. Complete the `feedForward()` function in `MLPNet.cpp`. You will need to write SIMD-ized `affineTransform`, `Softmax` and `MatMul2D1D` functions for this.
2. Implement a SIMD-ized ReLU activation as in case 1 of `applyActivation()` function in `MLPNet.cpp`. The ReLU activation function is defined as

$$ReLU(x) = \begin{cases} x & \text{if } x \geq 0, \\ 0 & \text{otherwise} \end{cases} . \quad (1)$$

When the argument to ReLU is a vector, it is applied element-wise. Also write the back-propagation for ReLU using

$$ReLU'(x) = \begin{cases} 1 & \text{if } x > 0, \\ 0 & \text{otherwise} \end{cases} . \quad (2)$$

3. Run the network for topology = {784, 10}, 10000 training examples, 10000 testing examples, batch size = 40, number of epochs = 10, activation function ReLU. Write down the peak training and testing accuracy.

	Part of the source code of P4_F32vec4.h
44	<code>class F32vec4</code>
45	<code>{</code>
46	<code>public:</code>
47	
48	<code>__m128 v;</code>
49	
50	<code>float & operator[](int i){ return (reinterpret_cast<float*>(&v))[i]; }</code>
51	<code>float operator[](int i) const { return (reinterpret_cast<const</code>
52	<code>float*>(&v))[i]; }</code>
53	<code>F32vec4():v(_mm_set_ps1(0)){} F32vec4(const __m128 &a):v(a) {} F32vec4(const float &a):v(_mm_set_ps1(a)) {}</code>
54	
55	
56	
57	<code>F32vec4(const float &f0, const float &f1, const float &f2, const float</code>
...	<code>&f3):v(_mm_set_ps(f3,f2,f1,f0)) {}</code>
62	<code>/* Arithmetic Operators */</code>
63	<code>friend F32vec4 operator +(const F32vec4 &a, const F32vec4 &b) { return</code>
...	<code>_mm_add_ps(a,b); }</code>
72	<code>/* Square Root */</code>
73	<code>friend F32vec4 sqrt (const F32vec4 &a){ return _mm_sqrt_ps (a); }</code>
...	
214	<code>} __attribute__ ((aligned(16)));</code>
215	
216	
217	<code>typedef F32vec4 fvec;</code>
218	<code>typedef float fscal;</code>
219	<code>const int fvecLen = 4;</code>

Headers with vector instructions emulated by scalar operations are provided in **PSEUDO_F32vec4.h** (vector with 4 entries) and **PSEUDO_F32vec1.h** (vector with 1 entry). They can be used in the similar way as **P4_F32vec4.h**, but won't give any speed up. They supposed to be used for debugging and comparison.

2_SIMD/0_Matrix: description

The Matrix exercise requires to parallelize square root extraction over a set of float variables arranged in a square matrix, see Fig. 1. The initial code gives an implemented scalar part and leaves blank space for a vector part. Therefore the initial output shows 0 time for vector calculations and infinite speed up factor, which should be currently ignored.

$$a = \begin{bmatrix} a_{00} & a_{01} & a_{02} \\ a_{10} & a_{11} & a_{12} \\ a_{20} & a_{21} & a_{22} \end{bmatrix} \quad c = \begin{bmatrix} \sqrt{a_{00}} & \sqrt{a_{01}} & \sqrt{a_{02}} \\ \sqrt{a_{10}} & \sqrt{a_{11}} & \sqrt{a_{12}} \\ \sqrt{a_{20}} & \sqrt{a_{21}} & \sqrt{a_{22}} \end{bmatrix} = ?$$

Fig. 1. The explanation of Matrix task.

	Part of the source code of Matrix.cpp
23	<code>float a[N][N]; // input array</code>
24	<code>float c[N][N]; // output array for scalar computations</code>
25	<code>float c_simd[N][N]; // output array for SIMD computations</code>

	Part of the source code of Matrix.cpp
26	
27	template<typename T> // required calculations
28	T f(T x) {
29	return sqrt(x);
30	}
...	
45	int main() {
46	
47	// fill classes by random numbers
48	for(int i = 0; i < N; i++) {
49	for(int j = 0; j < N; j++) {
50	a[i][j] = float(rand())/float(RAND_MAX); // put a random value, from 0 to
51	1
52	}
53	}
54	/// -- CALCULATE --
55	/// SCALAR
56	TStopwatch timerScalar;
57	for(int ii = 0; ii < NIter; ii++)
58	for(int i = 0; i < N; i++) {
59	for(int j = 0; j < N; j++) {
60	c[i][j] = f(a[i][j]);
61	}
62	}
63	timerScalar.Stop();
64	
65	/// SIMD VECTORS
66	TStopwatch timerSIMD;
67	// TODO
68	timerSIMD.Stop();
69	
70	double tScal = timerScalar.RealTime()*1000;
71	double tSIMD1 = timerSIMD.RealTime()*1000;
72	
73	cout << "Time scalar: " << tScal << " ms " << endl;
74	cout << "Time SIMD: " << tSIMD1 << " ms, speed up " << tScal/tSIMD1 <<
75	endl;
76	CheckResults(c,c_simd);
77	
78	return 1;
79	}
	Typical output
	Time scalar: 599.968 ms
	Time SIMD: 0 ms, speed up inf
	ERROR! SIMD and scalar results are not the same.

2_SIMD/0_Matrix: solution

For vector calculations all $N \times N$ scalar variables need to be divided into $N \times (N/4)$ groups with 4 elements, that will be treated as vectors. The solution groups scalars within each row, thus creating a matrix with N rows and $N/4$ columns.

Initial data copying into vectors can be avoided by reinterpreting them as vectors in case the scalar data is aligned to 16 bytes (line 23 of the solution). In the same way copying can be avoided at the output stage (line 25).

Although one can use the same algorithm structure with a triple loop as in the scalar code, we iterate over 4 elements groups rather than over each matrix element, therefore the step of the innermost loop need to be modified.

	Part of the source code of Matrix_solution.cpp
23	<code>float a[N][N] __attribute__((aligned(16))); // input array</code>
24	<code>float c[N][N]; // output array for scalar</code>
	<code>computations</code>
25	<code>float c_simd[N][N] __attribute__((aligned(16))); // output array for SIMD</code>
	<code>computations</code>
...	
66	<code>TStopwatch timerSIMD;</code>
67	<code>for(int ii = 0; ii < NIter; ii++)</code>
68	<code>for(int i = 0; i < N; i++) {</code>
69	<code>for(int j = 0; j < N; j+=fvecLen) {</code>
70	<code>fvec &aVec = reinterpret_cast<fvec>(a[i][j]);</code>
71	<code>fvec &cVec = reinterpret_cast<fvec>(c_simd[i][j]);</code>
72	<code>cVec = f(aVec);</code>
73	<code>}</code>
74	<code>}</code>
75	<code>timerSIMD.Stop();</code>
	Typical output after solution
	Time scalar: 604.324 ms
	Time SIMD: 150.4 ms, speed up 4.01811
	SIMD and scalar results are the same.

Since the function `f(...)`, which must be applied to the data, is a template, we can use the same function call format.

Since 4 float variables fit into a single SIMD vector, all calculations are done in parallel and no overhead operations are required, the expected speed-up factor should be 4.

2_SIMD/1_QuadraticEquation: description

The QuadraticEquation exercise requires to vectorize solution of a set of quadratic equations in four different ways: using either (1) copying of data or (2) casting and using either (3) SIMD intrinsics or (4) `fvec` type from the header file. Then (5) compare the calculation time depending on amount of processed data.

It is recommended to compile the code with `-fno-tree-vectorize` option to prevent auto-vectorization of the scalar code, otherwise comparison of the vectorized and scalar codes will not be direct.

	Part of the source code of QuadraticEquation.cpp
74	<code>// fill parameters by random numbers</code>
75	<code>for(int i = 0; i < N; i++) {</code>
76	<code>a[i] = float(rand())/float(RAND_MAX); // put a random value, from 0 to 1</code>

	Part of the source code of QuadraticEquation.cpp
77	b[i] = float(rand())/float(RAND_MAX);
78	c[i] = -float(rand())/float(RAND_MAX);
79	}
80	
81	/// -- CALCULATE --
82	
83	// scalar calculations
84	TStopwatch timerScalar;
85	for(int io=0; io<NIterOut; io++)
86	for(int i=0; i<N; i++)
87	{
88	float det = b[i]*b[i] - 4*a[i]*c[i];
89	x[i] = (-b[i]+sqrt(det))/(2*a[i]);
90	}
91	timerScalar.Stop();
	Typical output
	Time scalar: 431.702 ms
	Time SIMD1: 0 ms, speed up inf
	Time SIMD2: 0.000953674 ms, speed up 452672
	Time SIMD3: 0 ms, speed up inf
	Time SIMD4: 0.000953674 ms, speed up 452672
	ERROR! SIMD1 and scalar results are not the same.
	ERROR! SIMD2 and scalar results are not the same.
	ERROR! SIMD3 and scalar results are not the same.
	ERROR! SIMD4 and scalar results are not the same.

2_SIMD/1_QuadraticEquation: solution

(1) The first task is to vectorize calculations using copying of data and SIMD intrinsics.

One need to use `__m128` type and `_mm_set_ps` function to copy the data. This function takes four float arguments and fills 4 32-bit elements in `__m128` variable starting from the last one. Therefore for i-th vector one uses (i*fvecLen+3)-th, (i*fvecLen+2)-th, (i*fvecLen+1)-th and (i*fvecLen)-th scalars as arguments.

The `_mm_set_ps1` function is used to convert constants 4 and 2 into the SIMD type. The `_mm_sub_ps`, `_mm_mul_ps`, `_mm_sqrt_ps` and `_mm_div_ps` are used for all required calculations.

	Part of the source code of QuadraticEquation_solution.cpp
98	<code>__m128 aV = _mm_set_ps(a[i*fvecLen+3],a[i*fvecLen+2],a[i*fvecLen+1],a[i*fvecLen]);</code>
99	<code>__m128 bV = _mm_set_ps(b[i*fvecLen+3],b[i*fvecLen+2],b[i*fvecLen+1],b[i*fvecLen]);</code>
100	<code>__m128 cV = _mm_set_ps(c[i*fvecLen+3],c[i*fvecLen+2],c[i*fvecLen+1],c[i*fvecLen]);</code>
101	
102	<code>const __m128 det = _mm_sub_ps(_mm_mul_ps(bV,bV) ,</code> <code>_mm_mul_ps(_mm_set_ps1(4),_mm_mul_ps(aV,cV)));</code>
103	<code>__m128 xV =</code> <code>_mm_div_ps(_mm_sub_ps(_mm_sqrt_ps(det),bV),_mm_mul_ps(_mm_set_ps1(2),aV));</code>

(2) The second task is vectorisation using cast and SIMD intrinsics.

Since the data is already aligned to 16 bytes one can use `reinterpret_cast` directly to the part of the scalar input and output arrays starting from the i-th element.

The calculation part remains the same.

	Part of the source code of QuadraticEquation_solution.cpp
114	<code>for(int i=0; i<N; i+=fvecLen)</code>
115	<code>{</code>
116	<code>__m128& aV = (reinterpret_cast<__m128&>(a[i]));</code>
117	<code>__m128& bV = (reinterpret_cast<__m128&>(b[i]));</code>
118	<code>__m128& cV = (reinterpret_cast<__m128&>(c[i]));</code>
119	
120	<code>__m128& xV = (reinterpret_cast<__m128&>(x_simd2[i]));</code>

(3) The third task is to vectorize calculations using copying of data and headers.

With headers one uses **fvec** type and **fvec** constructor to copy the data. The constructor takes 4 float arguments and fills 4 32-bit elements starting from the first one. Therefore for the i-th vector one uses (i*fvecLen)-th, (i*fvecLen+1)-th, (i*fvecLen+2)-th and (i*fvecLen+3)-th scalars as arguments.

The calculations part remains the same.

	Part of the source code of QuadraticEquation_solution.cpp
134	<code>fvec aV = fvec(a[i*fvecLen],a[i*fvecLen+1],a[i*fvecLen+2],a[i*fvecLen</code>
135	<code>+3]); fvec bV = fvec(b[i*fvecLen],b[i*fvecLen+1],b[i*fvecLen+2],b[i*fvecLen</code>
136	<code>+3]); fvec cV = fvec(c[i*fvecLen],c[i*fvecLen+1],c[i*fvecLen+2],c[i*fvecLen</code>
137	<code>+3]);</code>
138	<code>const fvec det = bV*bV - 4*aV*cV;</code>
139	<code>fvec xV = (-bV+sqrt(det))/(2*aV);</code>

(4) The fourth task is to vectorize using cast of data and headers. It can be done as in the tasks (2) and (3).

	Typical output after solution
	Time scalar: 432.194 ms
	Time SIMD1: 109.914 ms, speed up 3.93211
	Time SIMD2: 105.318 ms, speed up 4.1037
	Time SIMD3: 111.117 ms, speed up 3.88954
	Time SIMD4: 105.588 ms, speed up 4.09321
	SIMD1 and scalar results are the same.
	SIMD2 and scalar results are the same.
	SIMD3 and scalar results are the same.
	SIMD4 and scalar results are the same.

It is expected to have exactly the same speed up factor independently of SSE intrinsics overloading. The speed up factor of 4 should be achieved with **reinterpret_cast**. It should be slightly less in the other case, because additional time is spend on data management.

(5) The initial number of vectors is 10,000, which gives the size of arrays 10000*16 bytes = 0.16 MB. All 5 arrays should fit into the cache memory (typical size 1-20 MB), therefore the speed-up factor of 4 should be achievable with **reinterpret_cast**.

When the number of vectors is increased by a factor of 1000, the size of arrays becomes 160 MB, then the arrays will be stored in the much slower RAM memory. Therefore the speed-up factor will be less than 4. The speed-up factor should remain the same (low) also for 1,000,000 vectors.

	Typical output after solution
	Time scalar: 4322.71 ms

	Typical output after solution
	Time SIMD1: 1197.12 ms, speed up 3.61093 Time SIMD2: 1130.53 ms, speed up 3.82361 Time SIMD3: 1215.68 ms, speed up 3.5558 Time SIMD4: 1137.61 ms, speed up 3.79981 SIMD1 and scalar results are the same. SIMD2 and scalar results are the same. SIMD3 and scalar results are the same. SIMD4 and scalar results are the same.

2_SIMD/2_CheckSum: description

The CheckSum exercise requires to vectorize the check sum calculation for the array of data. Check sum is defined as XOR-sum over all bytes of the array (lines 29-30 in the code below). In addition, consider if parallel calculations here can be done without the hardware support of parallelization.

	Part of the source code of CheckSum.cpp
23	<code>template< typename T ></code>
24	<code>T Sum(const T* data, const int N)</code>
25	<code>{</code>
26	<code> T sum = 0;</code>
27	<code></code>
28	<code> for (int i = 0; i < N; ++i)</code>
29	<code> sum = sum ^ data[i];</code>
30	<code> return sum;</code>
31	<code>}</code>
32	<code></code>
33	<code>int main() {</code>
34	<code></code>
35	<code> // fill string by random values</code>
36	<code> for(int i = 0; i < N; i++) {</code>
37	<code> str[i] = 256 * (double(rand()) / RAND_MAX); // put a random value, from 0</code>
38	<code>to 255</code>
39	<code> }</code>
40	<code> /// -- CALCULATE --</code>
41	<code></code>
42	<code> /// SCALAR</code>
43	<code></code>
44	<code> unsigned char sumS = 0;</code>
45	<code> TStopwatch timerScalar;</code>
46	<code> for(int ii = 0; ii < NIter; ii++)</code>
47	<code> sumS = Sum<unsigned char>(str, N);</code>
48	<code> timerScalar.Stop();</code>

2_SIMD/2_CheckSum: solution

The parallelization can be achieved, since XOR operator is applied bitwise and can be applied to the SIMD register data in the same way as to `char` variable, such that we can use `__m128` type to pack 16 `char` variables into a vector and treat them simultaneously.

The input **char** array should be reinterpreted as an array of **fvecs**. Then XOR sum of **fvecs** is calculated using the given template function. The resulting **fvec** variable should be reinterpreted back as an array of 16 **char** variables and the scalar sum of 16 elements is calculated directly.

	Part of the source code of CheckSum_solution.cpp
52	<code>unsigned char sumV = 0;</code>
53	
54	<code>const int fvecCharLen = fvecLen*4;</code>
55	<code>const int NV = N/fvecCharLen;</code>
56	
57	<code>TStopwatch timerSIMD;</code>
58	<code>for(int ii = 0; ii < NIter; ii++) {</code>
59	<code> fvec sumVV = 0;</code>
60	<code> sumVV = Sum<fvec>(reinterpret_cast<fvec*>(str), NV);</code>
61	<code> unsigned char *sumVS = reinterpret_cast<unsigned char*>(&sumVV);</code>
62	
63	<code> sumV = sumVS[0];</code>
64	<code> for (int iE = 1; iE < fvecCharLen; ++iE)</code>
65	<code> sumV ^= sumVS[iE];</code>
66	<code>}</code>
67	<code>timerSIMD.Stop();</code>

Parallelization without SIMD instructions can be achieved, since XOR operator is applied bitwise and can be applied to integer data in the same way as to SIMD data and byte data. The only difference is the length of the vector, that is 4 in case of integer.

	Part of the source code of CheckSum_solution.cpp
71	<code>unsigned char sumI = 0;</code>
72	
73	<code>const int intCharLen = 4;</code>
74	<code>const int NI = N/intCharLen;</code>
75	
76	<code>TStopwatch timerINT;</code>
77	<code>for(int ii = 0; ii < NIter; ii++) {</code>
78	<code> int sumII = Sum<int>(reinterpret_cast<int*>(str), NI);</code>
79	<code> unsigned char *sumIS = reinterpret_cast<unsigned char*>(&sumII);</code>
80	
81	<code> sumI = sumIS[0];</code>
82	<code> for (int iE = 1; iE < intCharLen; ++iE)</code>
83	<code> sumI ^= sumIS[iE];</code>
84	<code>}</code>
85	<code>timerINT.Stop();</code>
	Typical output after solution
	Time scalar: 274.435 ms
	Time INT: 65.3081 ms, speed up 4.20216
	Time SIMD: 25.3839 ms, speed up 10.8114
	Results are the same.

A typical speed-up with integer should be 4, as expected, since 4 bytes are packed into one integer and time for the additional loop with 3 iterations is negligible. The speed-up factor using **fvec** theoretically should be 16, practically the achieved speed-up factor is about 10 to 12.

	Typical output of KFLineFitter_solution2_simd.cpp
	Begin
	Time: 0.647068 ms
	End