

Note 17-11

First Steps with OpenMP 4.5 on Ubuntu and Nvidia GPUs

Stefan Rosenberger

May 16, 2018

CONTENTS

1	Setup System	2
1.1	Install OpenMP for GPU's	2
1.2	Run a test example	4
2	Test OpenMP for the sandbox example	6
2.1	Working Test CODE	7
3	CRS Matrix-Vector mutliplication	8
3.1	Compiler options	11
4	Diagonal-Matrix Vector mutliplication	13
5	ELLPACK Matrix Vector mutliplication	14
6	Scalar Product	17
6.1	Linear Clause	18

Preface: We want to emphasize that this document is a note on our OpenMP 4.5 parallelization. Therefore, we check only *sloppy* with respect to spellings and formulations.

1 SETUP SYSTEM

We use the following system to install and run OpenMP 4.5:

- OS:

```
1 No LSB modules are available.
2 Distributor ID: LinuxMint
3 Description: Linux Mint 18.1 Serena
4 Release: 18.1
5 Codename: serena
```

- Host:

```
1 /0/3d          memory      32GiB System Memory
2 /0/3d/3        memory      16GiB DIMM Synchronous 2133 MHz (0.5 ns)
3 /0/43          memory      256KiB L1 cache
4 /0/44          memory      1MiB L2 cache
5 /0/45          memory      8MiB L3 cache
6 /0/46          processor   Intel(R) Core(TM) i7-7700 CPU @ 3.60GHz
7 /0/1/0.0.0     /dev/sda  disk        250GB Samsung SSD 850
```

- GPU:

```
1 Mon Sep  4 14:59:44 2017
2 +-----+
3 | NVIDIA-SMI 375.82                Driver Version: 375.82                |
4 +-----+-----+-----+-----+
5 | GPU   Name           Persistence-M| Bus-Id        Disp.A | Volatile Uncorr. ECC |
6 | Fan  Temp  Perf    Pwr:Usage/Cap|      Memory-Usage | GPU-Util  Compute M. |
7 +-----+-----+-----+-----+
8 |    0  GeForce GTX 106...    Off | 0000:01:00.0    On   |           N/A       |
9 | 40%   34C    P0      28W / 120W | 499MiB / 6072MiB |      3%      Default |
10 +-----+-----+-----+-----+
11
12 +-----+
13 | Processes:                                             GPU Memory |
14 |  GPU       PID    Type   Process name                               Usage      |
15 +-----+-----+-----+-----+
16 |    0       1491    G      /usr/lib/xorg/Xorg                           286MiB |
17 |    0       1994    G      cinnamon                                    206MiB |
18 |    0       3605    G      /usr/lib/firefox/firefox                      1MiB |
19 +-----+-----+-----+-----+
```

1.1 INSTALL OPENMP FOR GPU'S

First of all one needs to setup the local system for OpenMP 4.5. Change gcc version:

```
1 sudo add-apt-repository ppa:ubuntu-toolchain-r/test
2 sudo apt update
3 sudo apt install gcc-6
```

```
4
5 sudo update-alternatives --install /usr/bin/gcc gcc /usr/bin/gcc-6 60 --slave /usr/bin/g++
6      g++ /usr/bin/g++-6
```

Furthermore we use the bash script (from [3])

```
1  #!/bin/sh
2
3  #
4  # Build GCC with support for offloading to NVIDIA GPUs.
5  #
6
7  work_dir=$HOME/offload/wrk
8  install_dir=$work_dir/install
9
10 # Location of the installed CUDA toolkit
11 cuda=/usr/local/cuda
12
13 # Build assembler and linking tools
14 mkdir -p $work_dir
15 cd $work_dir
16 git clone https://github.com/MentorEmbedded/nvptx-tools
17 cd nvptx-tools
18 ./configure \
19     --with-cuda-driver-include=$cuda/include \
20     --with-cuda-driver-lib=$cuda/lib64 \
21     --prefix=$install_dir
22 make
23 make install
24 cd ..
25
26 # Set up the GCC source tree
27 git clone https://github.com/MentorEmbedded/nvptx-newlib
28 svn co svn://gcc.gnu.org/svn/gcc/trunk gcc
29 cd gcc
30 contrib/download_prerequisites
31 ln -s ../nvptx-newlib/newlib newlib
32 cd ..
33 target=$(gcc/config.guess)
34
35 # Build nvptx GCC
36 mkdir build-nvptx-gcc
37 cd build-nvptx-gcc
38 ../gcc/configure \
39     --target=nvptx-none --with-build-time-tools=$install_dir/nvptx-none/bin \
40     --enable-as-accelerator-for=$target \
41     --disable-sjlj-exceptions \
42     --enable-newlib-io-long-long \
43     --enable-languages="c,c++,fortran,lto" \
44     --prefix=$install_dir
45 make -j4
46 make install
47 cd ..
48
49 # Build host GCC
50 mkdir build-host-gcc
51 cd build-host-gcc
52 ../gcc/configure \
53     --enable-offload-targets=nvptx-none \
```

```

54  --with-cuda-driver-include=$cuda/include \
55  --with-cuda-driver-lib=$cuda/lib64 \
56  --disable-bootstrap \
57  --disable-multilib \
58  --enable-languages="c,c++,fortran,lto" \
59  --prefix=$install_dir
60  make -j4
61  make install
62  cd ..

```

1.2 RUN A TEST EXAMPLE

We test the following source code:

```

1  #include <set>
2  #include <iostream>
3  #include <assert.h>
4  #include <vector>
5
6  using namespace std;
7
8  void vec_mult(int N, const int n_print=10)
9  {
10     double p[N], v1[N], v2[N];
11
12     for(int ii=0; ii<N; ++ii){
13         p[ii] = ii%5;
14         v1[ii] = ii%6;
15         v2[ii] = ii%7;
16     }
17
18     int i;
19     #pragma omp target map(to: v1[0:N], v2[0:N], p[0:N])
20     #pragma omp parallel for private(i)
21     for (i=0; i<N; i++){
22         p[i] = v1[i] * v2[i];
23     }
24
25     int num_print = 0;
26     if(n_print > N) num_print = N;
27     else num_print = n_print;
28     for(int ii=0; ii<num_print; ++ii) cout << p[ii] << "  ";
29     cout << endl;
30 }
31
32
33 /*
34  * ===== MAIN =====
35  * Test OpenMP 4.5
36  */
37 int main(int argc, char **argv)
38 {
39     cout << "#####" << endl;
40     cout << "#####Test_Start#####" << endl;
41     cout << "#####" << endl << endl;
42
43
44     vec_mult(100000);

```

```

45
46     cout << "Test complete!\n";
47     return 0;
48 }

```

We compile and run the code with the makefile:

```

1  default:
2      rm -f a.out
3      /home/rosenbs/offload/wrk/install/bin/g++ -std=c++11 -O3 -fopenmp -DOPENMP
4          -foffload=nvptx-none -Wall test.cpp
5      ./a.out
6
7  run:
8      ./a.out
9
10 clean:
11     rm -f a.out

```

Note, in line 3 one **has to use the g++ compiler** which is created from the bash file.

We get the runtime output:

```

1  ./a.out
2  #####
3  #####      Test Start      #####
4  #####
5
6  0  1  2  3  4  0  1  2  3  4
7  Test complete!

```

To check that we are really on the GPU we test the code with nvprof:

```

1  rosenbs@math068 ~/src/test/testomp $ nvprof ./a.out
2  #####
3  #####      Test Start      #####
4  #####
5
6  ==14296== NVPROF is profiling process 14296, command: ./a.out
7  0  1  2  3  4  0  1  2  3  4
8  Test complete!
9  ==14296== Profiling application: ./a.out
10 ==14296== Profiling result:
11 Time(%)      Time      Calls      Avg      Min      Max  Name
12  95.25%    3.8015ms        1  3.8015ms  3.8015ms  3.8015ms  _Z8vec_multii$_omp_fn$0
13   4.75%    189.63us         7  27.090us   544ns  62.655us  [CUDA memcpy HtoD]
14
15 ==14296== API calls:
16 Time(%)      Time      Calls      Avg      Min      Max  Name
17  59.93%    104.82ms         1  104.82ms  104.82ms  104.82ms  cuCtxCreate
18  34.01%    59.484ms         1  59.484ms  59.484ms  59.484ms  cuCtxDestroy
19   2.20%     3.8485ms         1  3.8485ms  3.8485ms  3.8485ms  cuCtxSynchronize
20   1.20%     2.1007ms        22  95.485us  15.600us  886.99us  cuLinkAddData
21   0.92%     1.6123ms         1  1.6123ms  1.6123ms  1.6123ms  cuModuleLoadData
22   0.42%     729.89us         1  729.89us  729.89us  729.89us  cuLinkComplete
23   0.26%     457.66us         1  457.66us  457.66us  457.66us  cuMemAllocHost
24   0.25%     430.49us         1  430.49us  430.49us  430.49us  cuLaunchKernel
25   0.25%     428.78us         7  61.254us  4.1860us  133.51us  cuMemcpyHtoD
26   0.24%     427.22us         1  427.22us  427.22us  427.22us  cuLinkCreate
27   0.15%     256.86us         1  256.86us  256.86us  256.86us  cuMemFreeHost
28   0.09%     149.88us         2  74.940us  5.0250us  144.85us  cuMemAlloc

```

29	0.04%	76.635us	15	5.1090us	101ns	70.902us	cuDeviceGetAttribute
30	0.04%	64.396us	2	32.198us	6.4940us	57.902us	cuMemFree
31	0.00%	5.5170us	4	1.3790us	440ns	3.7350us	cuDeviceGetCount
32	0.00%	4.1140us	3	1.3710us	676ns	2.4450us	cuDeviceGet
33	0.00%	2.7490us	12	229ns	106ns	892ns	cuCtxGetDevice
34	0.00%	2.4150us	8	301ns	195ns	542ns	cuMemGetAddressRange
35	0.00%	2.2930us	1	2.2930us	2.2930us	2.2930us	cuLinkDestroy
36	0.00%	1.4190us	1	1.4190us	1.4190us	1.4190us	cuInit
37	0.00%	1.3550us	1	1.3550us	1.3550us	1.3550us	cuModuleGetFunction
38	0.00%	961ns	1	961ns	961ns	961ns	cuCtxGetCurrent
39	0.00%	837ns	1	837ns	837ns	837ns	cuMemHostGetDevicePointer
40	0.00%	557ns	2	278ns	233ns	324ns	cuFuncGetAttribute
41	0.00%	386ns	1	386ns	386ns	386ns	cuModuleGetGlobal

Additional settings (motivated by: Adam Dziekonski)

👉 adding the following was crucial to make it work with gcc 9.0 (experimental)
- add LIBDIR to the 'LD_LIBRARY_PATH' environment variable during execution
- add LIBDIR to the 'LD_RUN_PATH' environment variable during linking

But:

👉 This is NOT working for MPI applications!

2 TEST OPENMP FOR THE SANDBOX EXAMPLE

To let OpenMP run with MPI I need to change the default g++ compiler.

OLD:

```
1 g++ -> /etc/alternatives/g++*
```

NEW:

```
1 g++ -> /home/rosenbs/offload/wrk/install/bin/g++
```

Furthermore, I've created softlinks to all *wrapper* we have created in the first section, and put them into */usr/bin/*:

```
1 c++ -> /home/rosenbs/offload/wrk/install/bin/c++*
2 cpp -> /home/rosenbs/offload/wrk/install/bin/cpp*
3
4 g++ -> /home/rosenbs/offload/wrk/install/bin/g++
5 gcc -> /home/rosenbs/offload/wrk/install/bin/gcc
6
7 gcc-ar -> /home/rosenbs/offload/wrk/install/bin/gcc-ar*
8 gcc-nm -> /home/rosenbs/offload/wrk/install/bin/gcc-nm*
9 gcc-ranlib -> /home/rosenbs/offload/wrk/install/bin/gcc-ranlib*
10
11 gcov -> /home/rosenbs/offload/wrk/install/bin/gcov*
12 gcov-dump -> /home/rosenbs/offload/wrk/install/bin/gcov-dump*
13 gcov-tool -> /home/rosenbs/offload/wrk/install/bin/gcov-tool*
14
15 gfortran -> /home/rosenbs/offload/wrk/install/bin/gfortran
16
17 nvptx-none-ar -> /home/rosenbs/offload/wrk/install/bin/nvptx-none-ar*
18 nvptx-none-as -> /home/rosenbs/offload/wrk/install/bin/nvptx-none-as*
```

```

19 nvptx-none-ld -> /home/rosenbs/offload/wrk/install/bin/nvptx-none-ld*
20 nvptx-none-ranlib -> /home/rosenbs/offload/wrk/install/bin/nvptx-none-ranlib*
21 nvptx-none-run -> /home/rosenbs/offload/wrk/install/bin/nvptx-none-run*
22 nvptx-none-run-single -> /home/rosenbs/offload/wrk/install/bin/nvptx-none-run-single*
23
24
25 x86_64-pc-linux-gnu-accel-nvptx-none-gcc -> /home/rosenbs/offload/wrk/install/bin/x86_64
26 -pc-linux-gnu-accel-nvptx-none-gcc*
27 x86_64-pc-linux-gnu-c++ -> /home/rosenbs/offload/wrk/install/bin/x86_64-pc-linux-gnu-c++
28 x86_64-pc-linux-gnu-g++ -> /home/rosenbs/offload/wrk/install/bin/x86_64-pc-linux-gnu-g++
29 x86_64-pc-linux-gnu-gcc -> /home/rosenbs/offload/wrk/install/bin/x86_64-pc-linux-gnu-gcc*
30 x86_64-pc-linux-gnu-gcc-8.0.0 -> /home/rosenbs/offload/wrk/install/bin/x86_64-pc-linux-gnu-
31 gcc-8.0.0*
32 x86_64-pc-linux-gnu-gcc-ar -> /home/rosenbs/offload/wrk/install/bin/x86_64-pc-linux-gnu-gcc-ar*
33 x86_64-pc-linux-gnu-gcc-ranlib -> /home/rosenbs/offload/wrk/install/bin/x86_64-pc-linux-gnu-gcc-ranlib*
34 x86_64-pc-linux-gnu-gcc-ranlib -> /home/rosenbs/offload/wrk/install/bin/x86_64-pc-linux-gnu-gcc-ranlib*
35 x86_64-pc-linux-gnu-gfortran -> /home/rosenbs/offload/wrk/install/bin/x86_64-pc-linux-gnu-gfortran*

```

🔍 This is still NOT working for the sandbox example we have used for the paper *Effective OpenACC Parallelizaion for Sparce Matrix Problems!*
 Durning Link-step we get the error message:

```

1 x86_64-pc-linux-gnu-accel-nvptx-none-gcc: error: unrecognized command line option
2      '-nvptx-none'

```

2.1 WORKING TEST CODE

We create a *standalone* test example, which works for our MPI parallized solver.

Makefile:

```

1 default:
2     rm -f a.out
3     mpicxx -std=c++11 -O3 -fopenmp -DOPENMP -DOPENMP_OFFLOAD -foffload=nvptx-none test.cpp
4     ./a.out /home/rosenbs/src/test/Sandbox_Maxim/genBi_ellmat.bin /home/rosenbs/src/
5         test/Sandbox_Maxim/genBi_rhs_Ie.bin
6
7 run:
8     ./a.out /home/rosenbs/src/test/Sandbox_Maxim/genBi_ellmat.bin /home/rosenbs/src/
9         test/Sandbox_Maxim/genBi_rhs_Ie.bin
10
11 clean:
12     rm -f a.out

```

And in the file test.cpp we include our libraries:

```

1 include <set>
2 #include <iostream>
3 #include <assert.h>
4 #include <vector>
5
6
7 #include<stdlib.h>
8 #include<errno.h>
9 #include </home/rosenbs/src/sandbox_hackatron/PT_C/toolbox/toolbox_sandbox.h>
10 #include </home/rosenbs/src/sandbox_hackatron/PT_C/toolbox_funcs.hpp>

```

```

11
12 using namespace sr_solver;
13
14 #define LOAD_RHS
15 ....

```

The rest is copy-past of the my sanbox_example.cpp file.

☞ This works!

3 CRS MATRIX-VECTOR MUTLIPLICATION

We tried several different types of the Matrix-Vector parallelization with OpenMP 4.5

- First of all, we try the simplest:

```

1 #pragma omp target teams distribute parallel for
2 for (int i=0; i<sv; i++) {
3     const int *const __restrict p_col = col+dsp[i];
4     const double *const __restrict p_ele = ele+dsp[i];
5     double s = 0.0;
6     for (int j=0; j<dsp[i+1]-dsp[i]; j++) {
7         s += p_ele[j]*u[p_col[j]];
8     }
9     v[i] = s;
10 }

```

☞ – Solve time: 775ms

- With simd instruction

```

1 #pragma omp target teams distribute parallel for simd
2 for (int i=0; i<sv; i++) {
3     const int *const __restrict p_col = col+dsp[i];
4     const double *const __restrict p_ele = ele+dsp[i];
5     double s = 0.0;
6     for (int j=0; j<dsp[i+1]-dsp[i]; j++) {
7         s += p_ele[j]*u[p_col[j]];
8     }
9     v[i] = s;
10 }

```


☞ – Solve time: 9.3ms

- With simd instruction and reduction

```


1 #pragma omp target teams distribute parallel for simd
2 for (int i=0; i<sv; i++) {
3     const int *const __restrict p_col = col+dsp[i];
4     const double *const __restrict p_ele = ele+dsp[i];
5     double s = 0.0;
6     #pragma omp reduction(+:s)
7     for (int j=0; j<dsp[i+1]-dsp[i]; j++) {
8         s += p_ele[j]*u[p_col[j]];
9     }
10     v[i] = s;
11 }

```


 – Solve time: 9.01ms

- With ordered instruction

```
1 #pragma omp target teams distribute parallel for simd
2 for (int i=0; i<sv; i++) {
3     const int *const __restrict p_col = col+dsp[i];
4     const double *const __restrict p_ele = ele+dsp[i];
5     double s = 0.0;
6     #pragma omp ordered simd
7     for (int j=0; j<dsp[i+1]-dsp[i]; j++) {
8         s += p_ele[j]*u[p_col[j]];
9     }
10    v[i] = s;
11 }
```

 – Solve time: 35.7ms


- With simd in inner loop!

```
1 #pragma omp target teams distribute parallel for
2 for (int i=0; i<sv; i++) {
3     const int *const __restrict p_col = col+dsp[i];
4     const double *const __restrict p_ele = ele+dsp[i];
5     double s = 0.0
6     #pragma omp simd
7     for (int j=0; j<dsp[i+1]-dsp[i]; j++) {
8         s += p_ele[j]*u[p_col[j]];
9     }
10    v[i] = s;
11 }
```

 – Solve time: 8.84ms; WRONG RESULT!!! Race condition!

- With critical!

```
1 #pragma omp target teams distribute parallel for
2 for (int i=0; i<sv; i++) {
3     const int *const __restrict p_col = col+dsp[i];
4     const double *const __restrict p_ele = ele+dsp[i];
5     double s = 0.0
6     #pragma omp critical
7     for (int j=0; j<dsp[i+1]-dsp[i]; j++) {
8         s += p_ele[j]*u[p_col[j]];
9     }
10    v[i] = s;
11 }
```

 – *Endless loop*: break after 1 min of calculation! → a closer look at the loop — obvious!

- With reduction

```
1 #pragma omp target parallel for
2 for (int i=0; i<sv; i++) {
3     const int *const __restrict p_col = col+dsp[i];
4     const double *const __restrict p_ele = ele+dsp[i];
5     double s = 0.0
6     #pragma omp reduction(+:s)
```

```

7   for (int j=0; j<dsp[i+1]-dsp[i]; j++) {
8       s += p_ele[j]*u[p_col[j]];
9   }
10  v[i] = s;
11 }

```


 – Solve time: 114ms

- With simd instruction and reduction and num_teams definition

```

1  #pragma omp target teams distribute parallel for simd num_teams(kk)
2  for (int i=0; i<sv; i++) {
3      const int *const __restrict p_col = col+dsp[i];
4      const double *const __restrict p_ele = ele+dsp[i];
5      double s = 0.0
6      #pragma omp reduction(+:s)
7      for (int j=0; j<dsp[i+1]-dsp[i]; j++) {
8          s += p_ele[j]*u[p_col[j]];
9      }
10     v[i] = s;
11 }

```


 – Solve time kk=5: 10.4ms
– Solve time kk=10: 7.7ms
– Solve time kk=25: 8.3ms
– Solve time kk=50: 8.8ms
– Solve time kk=100: 8.8ms
– Solve time kk=200: 8.9ms
– Solve time kk=500: 8.9ms
– Solve time kk=1000: 8.9ms

- With simd instruction and reduction and num_teams, mun_threads definition

```

1  #pragma omp target teams distribute parallel for simd num_teams(10) mun_threads(kk)
2  for (int i=0; i<sv; i++) {
3      const int *const __restrict p_col = col+dsp[i];
4      const double *const __restrict p_ele = ele+dsp[i];
5      double s = 0.0
6      #pragma omp reduction(+:s)
7      for (int j=0; j<dsp[i+1]-dsp[i]; j++) {
8          s += p_ele[j]*u[p_col[j]];
9      }
10     v[i] = s;
11 }

```

 – Solve time kk=16: 7.6ms
– Solve time kk=32: 7.7ms
– Solve time kk=64: 7.6ms
– Solve time kk=128: 7.6ms
– Solve time kk=256: 7.6ms
– Solve time kk=512: 7.7ms

- Solve time kk=1024: 7.6ms
 - Solve time kk=2048: 7.8ms
- No influence on the parallelization!

- With simd instruction and reduction and num_teams, thread_limit definition

```

1 #pragma omp target teams distribute parallel for simd num_teams(10) thread_limit(kk)
2 for (int i=0; i<sv; i++) {
3     const int *const __restrict p_col = col+dsp[i];
4     const double *const __restrict p_ele = ele+dsp[i];
5     double s = 0.0
6 #pragma omp reduction(+:s)
7     for (int j=0; j<dsp[i+1]-dsp[i]; j++) {
8         s += p_ele[j]*u[p_col[j]];
9     }
10    v[i] = s;
11 }

```

- ☞ - Solve time kk=16: 7.7ms
- Solve time kk=32: 9.0ms
- Solve time kk=64: 9.0ms
- Solve time kk=128: 8.9ms
- Solve time kk=256: 9.0ms
- Solve time kk=512: 9.0ms
- Solve time kk=1024: 9.0ms
- Solve time kk=2048: 9.1ms

Only negative influence on the parallelization!

3.1 COMPILER OPTIONS

We consider the following kernel:

```

1 #pragma omp target teams distribute parallel for simd num_teams(10)
2 for (int i=0; i<sv; i++) {
3     const int *const __restrict p_col = col+dsp[i];
4     const double *const __restrict p_ele = ele+dsp[i];
5     double s = 0.0
6 #pragma omp reduction(+:s)
7     for (int j=0; j<dsp[i+1]-dsp[i]; j++) {
8         s += p_ele[j]*u[p_col[j]];
9     }
10    v[i] = s;
11 }

```


- Compile line:

mpicxx -std=c++11 -O3 -DNDEBUG -fopenmp -DOPENMP -foffload=nvptx-none test.cpp

- ☞ - Solve time: 7.7ms

- Add to compile line:

-fopenmp-simd

 – Solve time: 7.8ms

- Add to compile line:

-march=native

 – Solve time: 7.8ms


- Add to compile line (and replace O3):

-fopenmp-simd -O2 -march=native

 – Solve time: 7.7ms

- Add to compile line:

-mlong-double-64

 – Solve time: 7.65ms


- Add to compile line:

-m64

 – Solve time: 7.7ms


- Add to compile line:

-m64 -mlong-double-64

 – Solve time: 7.6ms


- Add to pragma num_threads (128) and add to compile line:

-mprefer-avx128

 – Solve time: 7.56ms


- Add to compile line:

-ftree-vectorize

 – Solve time: 7.61ms


- Add to compile line:

-ftree-vectorizer-verbose=1

 – Solve time: 7.61ms

- Add to compile line:

-m64 -mlong-double-64 -mprefer-avx128 -ftree-vectorize

 – Solve time: 7.61ms

4 DIAGONAL-MATRIX VECTOR MUTLIPLICATION


We tried several different types of the Diagonal-Matrix Vector parallelization with OpenMP 4.5. Wherein we calculate

$$v = \omega D u \quad (4.1)$$

wherein ω is a scalar.


- First of all, we try the simplest:

```
1 #pragma omp target parallel for
2 for (int ii = 0; ii < vs; ii++) {
3     v[ii] = omega * D[ii] * u[ii];
4 }
```

 – Solve time: 29.88ms


- With teams distribute

```
1 #pragma omp target teams distribute parallel for
2 for (int ii = 0; ii < vs; ii++) {
3     v[ii] = omega * D[ii] * u[ii];
4 }
```

 – Solve time: 1.97ms

- With simd

```
1 #pragma omp target teams distribute parallel for simd
2 for (int ii = 0; ii < vs; ii++) {
3     v[ii] = omega * D[ii] * u[ii];
4 }
```

 – Solve time: 0.86ms

- With num_teams

```
1 #pragma omp target teams distribute parallel for simd num_teams(kk)
2 for (int ii = 0; ii < vs; ii++) {
3     v[ii] = omega * D[ii] * u[ii];
4 }
```

- 🔍 – Solve time kk=5: 0.84ms
- Solve time kk=10: 0.72ms
- Solve time kk=25: 0.75ms
- Solve time kk=50: 0.77ms
- Solve time kk=100: 0.76ms
- Solve time kk=200: 0.77ms

- With num_threads

```
1 #pragma omp target teams distribute parallel for simd num_teams(10) num_threads(kk)
2 for (int ii = 0; ii < vs; ii++) {
3     v[ii] = omega * D[ii] * u[ii];
4 }
```

- 🔍 – Solve time kk=16: 0.72ms
- Solve time kk=32: 0.72ms
- Solve time kk=64: 0.72ms
- Solve time kk=128: 0.69ms
- Solve time kk=256: 0.71ms
- Solve time kk=512: 0.70ms
- Solve time kk=512: 0.70ms

We find the same setting as for the matrix vector product!

5 ELLPACK MATRIX VECTOR MUTLIPLICATION

We calculate the matrix vector product for ELLPACK type matrices.

- First of all, we try the simplest:

```
1 #pragma omp target teams distribute parallel for
2 for (int ii = 0; ii < v_size; ii += _block_size) {
3
4     const int stride = ii * max_length;
5     const int kk_max = (_block_size + ii < v_size) ? (_block_size + ii) : v_size;
6
7     for (int kk = ii; kk < kk_max; ++kk) {
8         S s = 0.0;
9
10        for (int jj = 0; jj < max_length; ++jj) {
11            const int index = stride + kk - ii + jj * _block_size;
12            s += ele[index] * u[col[index]];
13        }
14        v[kk] = s;
15    }
16 }
```


- 🔍 – Solve time: 87.67ms

- With the directives from the previous sections:

```

1  #pragma omp target teams distribute parallel for simd default(none) num_teams(10)
2      thread_limit(128)
3  for (int ii = 0; ii < v_size; ii += _block_size) {
4
5      const int stride = ii * max_length;
6      const int kk_max = (_block_size+ii < v_size) ? (_block_size+ii):v_size;
7
8      for (int kk = ii; kk < kk_max; ++kk) {
9          S s = 0.0;
10
11         for (int jj = 0; jj < max_length; ++jj) {
12             const int index = stride + kk - ii + jj * _block_size;
13             s += ele[index] * u[col[index]];
14         }
15         v[kk] = s;
16     }
17 }

```


 – Solve time: 47.65ms

- With reduction:

```

1  #pragma omp target teams distribute parallel for simd default(none) num_teams(10)
2      thread_limit(128)
3  for (int ii = 0; ii < v_size; ii += _block_size) {
4
5      const int stride = ii * max_length;
6      const int kk_max = (_block_size+ii < v_size)?(_block_size+ii):v_size;
7
8      for (int kk = ii; kk < kk_max; ++kk) {
9          S s = 0.0;
10
11         #pragma omp reduction(+:s)
12         for (int jj = 0; jj < max_length; ++jj) {
13             const int index = stride + kk - ii + jj * _block_size;
14             s += ele[index] * u[col[index]];
15         }
16         v[kk] = s;
17     }
18 }

```


 – Solve time: 46.16ms

- Second initialization:

```


1  #pragma omp target teams distribute default(none) num_teams(100) thread_limit(128)
2  for (int ii = 0; ii < v_size; ii += _block_size) {
3
4      const int stride = ii * max_length;
5      const int kk_max = (_block_size+ii < v_size)?(_block_size+ii):v_size;
6
7      #pragma omp parallel for simd
8      for (int kk = ii; kk < kk_max; ++kk) {
9          S s = 0.0;
10
11         #pragma omp reduction(+:s)
12         for (int jj = 0; jj < max_length; ++jj) {
13             const int index = stride + kk - ii + jj * _block_size;
14             s += ele[index] * u[col[index]];
15         }
16         v[kk] = s;
17     }
18 }

```

 – Solve time: 962.57ms


- Second initialization and define threads interior:

```
1 #pragma omp target teams distribute default(none) num_teams(10000)
2 for (int ii = 0; ii < v_size; ii += _block_size) {
3
4     const int stride = ii * max_length;
5     const int kk_max = (_block_size+ii < v_size)?(_block_size+ii):v_size;
6
7     #pragma omp parallel for num_threads(128)
8     for (int kk = ii; kk < kk_max; ++kk) {
9         S s = 0.0;
10
11     #pragma omp reduction(+:s)
12         for (int jj = 0; jj < max_length; ++jj) {
13             const int index = stride + kk - ii + jj * _block_size;
14             s += ele[index] * u[col[index]];
15         }
16         v[kk] = s;
17     }
18 }
```

 – Solve time: 87.47ms


- Second initialization with inner simd:

```
1 #pragma omp target teams distribute default(none) num_teams(10000)
2 for (int ii = 0; ii < v_size; ii += _block_size) {
3
4     const int stride = ii * max_length;
5     const int kk_max = (_block_size+ii < v_size)?(_block_size+ii):v_size;
6
7     #pragma omp parallel for simd num_threads(128)
8     for (int kk = ii; kk < kk_max; ++kk) {
9         S s = 0.0;
10
11     #pragma omp reduction(+:s)
12         for (int jj = 0; jj < max_length; ++jj) {
13             const int index = stride + kk - ii + jj * _block_size;
14             s += ele[index] * u[col[index]];
15         }
16         v[kk] = s;
17     }
18 }
```

 – Solve time: 27.15ms


- Change inner loops

```
1 _v.zero_acc();
2
3 #pragma omp target teams distribute default(none) num_teams(10000)
4 for (int ii = 0; ii < v_size; ii += _block_size) {
5     const int stride = ii * max_length;
6     const int kk_max = (_block_size+ii < v_size) ? (_block_size+ii):v_size;
7     for (int jj = 0; jj < max_length; ++jj) {
8         #pragma omp simd
9         for (int kk = ii; kk < kk_max; ++kk) {
10             const int index = stride + kk - ii + jj * _block_size;
11             v[kk] += ele[index] * u[col[index]];
12         }
13     }
14 }
```


 – Solve time: 32.63ms

- Define all 3 parallelizations

```
1 _v.zero_acc();
2
3 const int num_teams_local = v_size/_block_size + 1;
4 #pragma omp target teams distribute default(none) num_teams(60000)
5 for (int ii = 0; ii < v_size; ii += _block_size) {
6     const int stride = ii * max_length;
7     const int kk_max = (_block_size < v_size - ii) ?
8         (_block_size) : v_size - ii;
9     #pragma omp parallel for simd
10    for (int jj = 0; jj < max_length; ++jj) {
11        const int index = stride + jj * _block_size;
12        for (int kk = 0; kk < kk_max; ++kk) {
13            v[kk + ii] += ele[index + kk] * u[col[index + kk]];
14        }
15    }
16 }
```

 – Solve time: 41.26ms
Race condition!!!!

6 SCALAR PRODUCT

We try to improve our parallelization. The working version is

```
1 template<class S>
2 void scalar_product(const toolbox_vector<S> &x, const toolbox_vector<S> &y, S &s)
3 {
4     S s = 0.0;
5     const S *const __restrict x = _x.data(), *const __restrict y = _y.data();
6     const int x_size = _x.size();
7
8     #pragma omp target teams distribute parallel for simd map(tofrom: s)
9         reduction(+:s) schedule(simd:static, _block_size)
10    for(int i = 0; i < x_size; i++)
11    {
12        s += x[i] * y[i];
13    }
14    _s = s;
15 }
```

Since we observe with OpenACC 2 kernels for the execution, we try to create a *copy version* of a CUDA scalar product. Which means, we create a parallel executable region, and then a reduction step:

```
1 template<class S>
2 void scalar_product(const toolbox_vector<S> &x, const toolbox_vector<S> &y, S &s)
3 {
4     S s = 0.0;
5     const S *const __restrict x = _x.data(), *const __restrict y = _y.data();
6     const int x_size = _x.size();
7
8     const double* x_ptr = &x[0];
9     const double* y_ptr = &y[0];
10
11     const int loc_size = 1;
12
13     double loc_val[loc_size];
14     for(int ii=0; ii<loc_size; ++ii) loc_val[ii] = 0.0;
15     double* loc_ptr = &loc_val[0];
```

```

16
17 #pragma omp target enter data map(to:loc_ptr[0:loc_size])
18
19 #pragma omp target parallel for simd safelen(loc_size)
20     for(int ii=0; ii<x_size; ++ii){
21         const int jj = ii%loc_size;
22         loc_ptr[jj] += x_ptr[ii] * y_ptr[ii];
23     }
24 #pragma omp target parallel for reduction(+:s) map(tofrom:s)
25     for(int ii=0; ii<loc_size; ++ii) s += loc_ptr[ii];
26
27     _s = s;
28 }

```

☞ This version has a race condition. Note every try to use an atomic operation to fix the bug fails → compiling error.

An alternative can be a use of *number of threads*, therefore we try to call on the target in a parallel region.

6.1 LINEAR CLAUSE

We tried:

```
1 omp_get_num_threads();
```

This leads to the bug:

☞ error: there are no arguments to 'omp_get_num_threads' that depend on a template parameter, so a declaration of 'omp_get_num_threads' must be available [-fpermissive]
nthreads = omp_get_num_threads();

As a consequence we try to use the *linear* clause. We found online the description:

simd construct – clause

linear(list[:linear-step])

- has **private** clause semantics, and also firstprivate and lastprivate semantics
- `int x=2;`
`#pragma omp simd linear(x:4)`
 for (int i=0; i<12; i++)
 ... = x;
 printf("%d\n", x);
- in each iteration, private x is initialized as $x = x_0 + i * 4$, where x_0 is the initial value of x before entering the SIMD construct
- the value of x in the sequentially last iteration is assigned to the original list item

```
#pragma omp simd linear(i:1)
for (i=0; i<N; i++)
    a[i] = b[i] * c[i];
```

i is incremented
by 1 in every
iteration

```
m = 1;
#pragma omp simd linear(m:2)
for (i=0; i<N; i++)
    a[i] = a[i]*m;
```

```
a[0] = a[0] * (1+0*2);
a[1] = a[1] * (1+1*2);
a[2] = a[2] * (1+1*3);
...
```

50

Figure 6.1: Linear clause for OpenMP4.0 (Picture from [1])

Therefore we try:

```
1  const int loc_size = 12;
2  double loc_val[loc_size];
3  int x = 1;
4
5  #pragma omp simd linear(x:4)
6  for(int i=0; i<loc_size; i++) loc_val[i] = static_cast<double>(x);
7
8  for(int ii=0; ii<loc_size; ++ii) cout << loc_val[ii] << "  ";
```

and get the output:

```
1 1 1 1 1 1 1 1 1 1 1 1
```

Therefore the description of 6.1 is wrong. Lets try something else:

```
1  const S *const __restrict x = _x.data(), *const __restrict y = _y.data();
2  const int x_size = _x.size();
3
4  const int loc_size = 12;
5  double loc_val[loc_size];
6  for(int ii=0; ii<loc_size; ++ii) loc_val[ii] = 0.0;
7  double* loc_ptr = &loc_val[0];
8
9  int m = 1;
10
11 #pragma omp target parallel for simd linear(m:2) map(tofrom:loc_ptr[0:loc_size])
12 for(int i=0; i<loc_size; i++){
13     loc_ptr[i] = x[i] + m;
14 }
15
16 for(int ii=0; ii<loc_size; ++ii) cout << loc_ptr[ii] << "  ";
```

wherein we initialize x with 1.0.

```
2 2 6 6 10 10 14 14 18 20 22 24
```

This is also not what I would expect from 6.1

We try (wherein we consider 20 elements of the array):

```
1  #pragma omp target parallel for simd linear(m:5) safelen(10) map(tofrom:loc_ptr[0:loc_size])
```

and get

```
2 2 2 17 17 17 32 32 32 47 47 47 62 62 72 72 82 82 92 92
```

By guessing I would conclude that the clause `linear` add in every iteration 5 to the value `m`, but not individual for every thread. It sums up the values for every thread until OpenMP starts a new stride.

```
Sh**: this means we do not have access to the thread ID of the parallelization ....
```

Possibly we can use it for the stride steps in ELLPACK

BUG fixing:

```
Out of history, #include<omp.h> was included as last in the headers ....
```

We tried:

```

1  #pragma omp target enter data map(to:loc_ptr[0:loc_size])
2  #pragma omp target parallel for simd map(tofrom: loc_size) num_threads(100)
3  for(int ii=0; ii<loc_size+10; ++ii){
4      loc_size = omp_get_num_threads();
5      const int nthread = omp_get_thread_num();
6      loc_ptr[nthread] = nthread;
7  }
8  cout << "val:░░" << loc_size << endl;
9
10 #pragma omp target update from(loc_ptr[0:loc_size])
11 for(int ii=0; ii<loc_size; ++ii) cout <<"░░" << loc_ptr[ii];

```

and get the output:

```

👉 val: 8
    0 0 0 0 0 0

```

This means we have a bug!

REFERENCES

- [1] Chapman, B., Eeachempati, D., Li, K.: OpenMP 4.0 features (2017), <http://extremecomputingtraining.anl.gov/files/2014/01/OpenMP-40-features-ATPESC-final-v2.pdf>, [Online; accessed 12-09-2017]
- [2] heise Developer: Accelerator offloading mit gcc (2017), <https://www.heise.de/developer/artikel/Accelerator-Offloading-mit-GCC-3317330.html>, [Online; accessed 29-08-2017]
- [3] Walfridsson, K.: Building gcc with support for nvidia ptx offloading (2017), <https://kristerw.blogspot.co.at/2017/04/building-gcc-with-support-for-nvidia.html>, [Online; accessed 28-08-2017]