

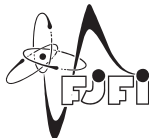
Programming GPU using TNL

Tomáš Oberhuber Jakub Klinkovský Radek Fučík
Aleš Wodecki

Department of Mathematics,
Faculty of Nuclear Sciences and Physical Engineering,
Czech Technical University in Prague



WSC 2019



Why GPU?



	Nvidia V100	Intel Xeon E5-4660
Cores	5120 @ 1.3GHz	16 @ 3.0GHz
Peak perf.	15.7/7.8 TFlops	0.4 / 0.2 TFlops
Max. RAM	32 GB	1.5 TB
Memory bw.	900 GB/s	68 GB/s
TDP	300 W	120 W

≈ 8,000 \$

Template Numerical Library

TNL = Template Numerical Library

- is written in C++ and profits from meta-programming
- provides unified interface to multi-core CPUs and GPUs (via CUDA)
- wants to be user friendly
- www.tnl-project.org
- \approx 300k lines of templated code
- documentation for few basic structures
- MIT license

Arrays

Arrays are basic structures for memory management

- `TNL::Array< ElementType, DeviceType, IndexType >`
- `DeviceType` says where the array resides
 - `TNL::Devices::Host` for CPU
 - `TNL::Devices::Cuda` for GPU
- memory allocation, I/O operations, elements manipulation ...

```
1 Array< float , Devices::Cuda, int > a( 100 );
2 auto lamda = [] __cuda_callable__ ( int elementIdx ) {
3     return elementIdx%5; }
4 a.evaluate( lambda );
```

Vectors

Vectors add algebraic operations to arrays:

- `TNL::Vector< RealType, DeviceType, IndexType >`
- addition, multiplication, scalar product, l_p norms ...

Vector and Array View

- arrays and vectors supports data sharing
- both are relatively complex structures
- TNL uses also lightweight counterparts `ArrayView`, `VectorView`
- both can be passed efficiently on GPU for example
- neither perform dynamic memory allocation/deallocation or deep copies

```
1 Vector< float , Devices::Cuda , int > v( 100 );  
2 VectorView< float , Devices::Cuda , int > view( v );
```

Parallel reduction

Parallel reduction is operation taking all array/vector elements as input and returns one value as output:

- array comparison
- scalar product
- l_p norm
- minimal/maximal value
- sum of all elements

```
1 float sum( 0.0 )
2 for( int i = 0; i < size; i++ )
3     sum += a[ i ];
```

Parallel reduction on GPU = 150 lines of code

```

1 // Parallel reduction on GPU
2 //
3 // This code is based on the code from the book "GPU Computing Gems"
4 // by John D. Owens, John L. Stone, and David B. Luebke.
5 //
6 // The code is licensed under the MIT license.
7 //
8 // Copyright (c) 2011, NVIDIA Corporation. All rights reserved.
9 //
10 // Permission is hereby granted, free of charge, to any person obtaining a copy
11 // of this software and associated documentation files (the "Software"), to deal
12 // in the Software without restriction, including without limitation the rights
13 // to use, copy, modify, merge, publish, distribute, sublicense, and/or sell
14 // copies of the Software, and to permit persons to whom the Software is
15 // furnished to do so, subject to the following conditions:
16 //
17 // The above copyright notice and this permission notice shall be included in
18 // all copies or substantial portions of the Software.
19 //
20 // THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS
21 // OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY,
22 // FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL
23 // THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER
24 // LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM,
25 // OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS
26 // IN THE SOFTWARE.
27 //
28 //
29 //
30 //
31 //
32 //
33 //
34 //
35 //
36 //
37 //
38 //
39 //
40 //
41 //
42 //
43 //
44 //
45 //
46 //
47 //
48 //
49 //
50 //
51 //
52 //
53 //
54 //
55 //
56 //
57 //
58 //
59 //
60 //
61 //
62 //
63 //
64 //
65 //
66 //
67 //
68 //
69 //
70 //
71 //
72 //
73 //
74 //
75 //
76 //
77 //
78 //
79 //
80 //
81 //
82 //
83 //
84 //
85 //
86 //
87 //
88 //
89 //
90 //
91 //
92 //
93 //
94 //
95 //
96 //
97 //
98 //
99 //
100 //
101 //
102 //
103 //
104 //
105 //
106 //
107 //
108 //
109 //
110 //
111 //
112 //
113 //
114 //
115 //
116 //
117 //
118 //
119 //
120 //
121 //
122 //
123 //
124 //
125 //
126 //
127 //
128 //
129 //
130 //
131 //
132 //
133 //
134 //
135 //
136 //
137 //
138 //
139 //
140 //
141 //
142 //
143 //
144 //
145 //
146 //
147 //
148 //
149 //
150 //

```


Parallel reduction in TNL

Take a look at scalar product:

```
1 float result( 0.0 );
2 for( int i = 0; i < size; i++ )
3     result += a[ i ] * b[ i ];
```

Let us rewrite it using C++ lambda functions as:

```
1 float a[ size ], b[ size ];
2
3 ...
4
5 auto fetch = [=] (int i)->float {
6     return a[i]*b[i]; };
7 auto reduce = [] (float& x, const float& y) {
8     x += y; };
9
10 float result( 0.0 );
11 for( int i = 0; i < size; i++ )
12     reduce( result , fetch( i ) );
```

Parallel reduction in TNL

Another example - l_p -norm:

```
1 const float p = 2.0;
2 float a[ size ];
3
4 auto fetch = [=] (int i)->float {
5     return pow( fabs( a[i] ), p ); };
6 auto reduce = [] (float& x, const float& y) {
7     x += y; };
8
9 float result( 0.0 );
10 for( int i = 0; i < size; i++ )
11     reduce( result , fetch( i ) );
```

Parallel reduction in TNL

Another example - arrays comparison:

```
1 bool zero = true;
2 const float p = 2.0;
3 float a[size], b[size];
4 ...
5 auto fetch = [=] (int i)->bool {
6     return ( a[i] == b[i] ); };
7 auto reduce = [] (float& x, const float& y) {
8     x = x && y; };
9
10 float result( zero );
11 for( int i = 0; i < size; i++ )
12     reduce( result , fetch( i ) );
```

Parallel reduction in TNL

To perform the same on GPU in TNL just add `__cuda_callable__` to lambdas...

```
1 auto fetch = [=] __cuda_callable__ (int i)->bool {  
2     return ( a[i] == b[i] ); };  
3 auto reduce = [] __cuda_callable__ (float& x,  
4                                     const float& y) {  
5     x = x && y; };
```

... and for certain reasons, deliver volatile version of reduce:

```
1 auto volatileReduce = [] __cuda_callable__ (  
2     volatile float& x,  
3     volatile const float& y) {  
4     x = x && y; };
```

This could be avoided when CUDA compiler supports C++17 better. Now call

```
1 Reduction< Devices::Cuda >::reduce( size , reduce ,  
2     volatileReduce , fetch , zero );
```

Expression Templates in TNL

Expression

$$\vec{x} = \vec{a} + 2\vec{b} + 3\vec{c}$$

can be evaluated in C as follows:

```
1 for( int i = 0; i < N; i++ )  
2     x[ i ] = a[ i ] + 2 * b[ i ] + 3 * c[ i ];
```

It is:

- efficient
- relatively simple
- works only on CPU - sequentially

Expression Templates in TNL

We can use operators overloading in C++:

```
1 x = a + 2 * b + 3 * c;
```

- it is very simple and easy to read
- can be performed in parallel on multicore CPUs or GPUs
- it is inefficient

Expression Templates in TNL

```
1 x = a + 2 * b + 3 * c;
```

The code is at the end performed almost like this:

```
1 Vector tmp1( N ), tmp2( N ), tmp3( N );
2 for( int i = 0; i < N; i++ )
3     tmp1[ i ] = 2 * b[ i ];
4 for( int i = 0; i < N; i++ )
5     tmp2[ i ] = 3 * c[ i ];
6 for( int i = 0; i < N; i++ )
7     tmp3[ i ] = tmp1[ i ] + tmp2[ i ];
8 for( int i = 0; i < N; i++ )
9     x[ i ] = a[ i ] + tmp3[ i ];
```

Expression Templates in TNL

We can use BLAS/cuBLAS:

```
1 cublasHandle_t handle;  
2 cublasSaxpy( handle, N, 1.0, a, 1, x, 1 );  
3 cublasSaxpy( handle, N, 2.0, b, 1, x, 1 );  
4 cublasSaxpy( handle, N, 3.0, c, 1, x, 1 );
```

- it is pretty hard to read
- works only for single precision
- more efficient than C++ version but still less efficient than C version

Expression Templates in TNL

The code is at the end performed almost like this:

```
1 for( int i = 0; i < tmp.size(); i++ )  
2     x[ i ] = a[ i ];  
3 for( int i = 0; i < tmp.size(); i++ )  
4     x[ i ] += 2 * b[ i ];  
5 for( int i = 0; i < tmp.size(); i++ )  
6     x[ i ] += 3 * c[ i ];
```

Expression Templates in TNL

	Memory	Operations	Conditions
C-like	$4N$	$5N$	N
C++	$7N$	$8N$	$4N$
Blas	$4N$	$5N$	$3N$

Expression Templates in TNL

Expression templates take the formula...

```
1 x = a + 2 * b + 3 * c;
```

... and parse it into a form of C++ typ:

```
1 Addition<  
2   Vector ,  
3   Addition<  
4     Multiplication< double , Vector > ,  
5     Multiplication< double , Vector > > > expr;
```

The expression is at the end evaluated like this:

```
1 for( int i = 0; i < N; i++ )  
2   x[ i ] = expr[ i ];
```

- it is simple and easy to read
- works for any type Real (float/double) and any Device (CPU/GPU)
- it is very efficient

Expression Templates in TNL

Expression templates take the formula...

```
1 x = a + 2 * b + 3 * c;
```

... and parse it into a form of C++ typ:

```
1 Addition<  
2   Vector ,  
3   Addition<  
4     Multiplication< double , Vector > ,  
5     Multiplication< double , Vector > > > expr;
```

The expression is at the end evaluated like this:

```
1 for( int i = 0; i < N; i++ )  
2   x[ i ] = expr[ i ];
```

- it is simple and easy to read
- works for any type Real (float/double) and any Device (CPU/GPU)
- it is very efficient

Expression Templates in TNL

	Memory	Operations	Conditions
C-like	$4N$	$5N$	N
C++	$7N$	$8N$	$4N$
Blas	$4N$	$5N$	$3N$
ET	$4N$	$5N$	N

Expression Templates & Parallel Reduction in TNL

Example:

```
1 using Vector = Vector< float , Devices::Cuda, int >;
2 using View = VectorView< float , Devices::Cuda, int >;
3 Vector av( 100 ), bv( 100 ), cv( 100 ), dv( 100 );
4 View a( av ), b( bv ), c( cv ), d( dv );
5 ...
6 float scalarProduct = ( a, b + 3 * c );
7 d = a + b * c + sin( d );
8 a = min( b, c );
9 float min_a = min( a );
10 float total_min = min( min( a, b ) );
```

Performance comparison

Performance was tested on:

- GPU Nvidia P100
 - 16 GB HBM2 @ 732 GB/s
 - 3584 CUDA cores, 4.7 TFlops in double precision
- CPU
 - AMD Ryzen 5 2600, 8MB L3 cache

Expression Templates in TNL

Scalar product: $r = (x, y)$.

Size	CPU			GPU		
	BLAS	TNL		cuBLAS	TNL	
	BW	BW	Speed-up	BW	BW	Speed-up
100k	20.8	2.7	0.1	49.3	69.9	1.41
200k	18.5	12.2	0.6	90.1	108.6	1.20
400k	18.4	13.3	0.7	142.2	159.1	1.11
800k	11.9	13.2	1.1	207.4	233.4	1.12
1.6M	13.6	15.6	1.1	313.6	333.3	1.06
3.2M	14.9	17.9	1.2	381.0	403.7	1.05
6.4M	16.7	17.0	1.0	417.1	431.8	1.03

Expression Templates in TNL

Vector addition: $x += a$.

Size	CPU			GPU		
	BLAS	TNL		cuBLAS	TNL	
	BW	BW	Speed-up	BW	BW	Speed-up
100k	46.0	11.6	0.2	152.2	174.8	1.14
200k	42.3	7.0	0.1	196.6	216.1	1.09
400k	16.9	29.3	1.7	277.6	294.4	1.06
800k	16.8	23.9	1.4	326.2	333.6	1.02
1.6M	17.3	25.1	1.4	362.5	374.2	1.03
3.2M	17.5	25.3	1.4	422.4	436.8	1.03
6.4M	17.4	25.7	1.4	456.6	469.8	1.02

Expression Templates in TNL

Vector addition: $x += a + b$.

Size	CPU			GPU		
	BLAS	TNL		cuBLAS	TNL	
	BW	BW	Speed-up	BW	BW	Speed-up
100k	30.3	29.3	0.9	188.3	190.9	1.01
200k	30.5	31.8	1.0	218.0	230.7	1.05
400k	13.7	32.8	2.3	243.1	305.9	1.25
800k	11.6	23.0	1.9	263.8	353.0	1.33
1.6M	11.7	24.4	2.0	285.9	389.4	1.36
3.2M	11.7	24.8	2.1	312.9	442.8	1.41
6.4M	11.7	25.7	2.1	327.3	471.9	1.44

Expression Templates in TNL

Vector addition: $x += a + b + c$.

Size	CPU			GPU		
	BLAS	TNL		cuBLAS	TNL	
	BW	BW	Speed-up	BW	BW	Speed-up
100k	25.4	7.7	0.30	194.7	236.5	1.21
200k	23.7	16.1	0.67	228.3	277.6	1.21
400k	13.0	31.0	2.38	218.3	330.9	1.51
800k	10.0	24.5	2.45	233.3	370.6	1.58
1.6M	9.9	23.6	2.38	249.6	403.4	1.61
3.2M	9.8	25.2	2.57	266.6	444.8	1.66
6.4M	9.8	25.9	2.64	276.6	471.3	1.70

More about TNL ...

TNL is available at

`www.tnl-project.org`

under MIT license.