

Programming GPU with TNL

Tomáš Oberhuber, J. Klinkovský, R. Fučík

Katedra matematiky,
Fakulta jaderná a fyzikálně inženýrská,
České vysoké učení technické v Praze



Winter Kindergarten School on LBM in Krakow 2022



Overview

CPU, GPU & LBM

TNL, Template Numerical Library

Memory management in TNL

Vector expressions

Parallel reduction

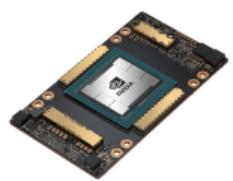
Solving the heat equation

Lattice Boltzmann Method

LBM - Lattice Boltzmann Method

- ▶ LBM is an efficient numerical method
 - ▶ it is easy to implement
 - ▶ perfectly profits from the GPU architecture

Comparison of GPU and CPU



	Nvidia A100	AMD Instinct MI100	AMD EPYC 2 Rome 7H12
Počet jader	6912 @ 1.41GHz	7680 @ 1.2GHz	64 @ 2.4GHz
Max. výkon	19.5/9.7 TFlops	23/11 TFlops	2 / 1 TFlops
Max. výkon s tenzory	156 /- TFlops	-/-	-/-
Max. RAM	80 GB	32 GB	2 TB
Datová propustnost	2039 GB/s	1200 GB/s	204 GB/s
Energetická náročnost	400 W	300 W	280 W

Unfortunately, programming GPU is difficult.

Template Numerical Library

TNL = Template Numerical Library

- ▶ numerical library for modern parallel architectures
- ▶ written in C++ and profiting from features of C++14 and C++17
- ▶ offers unified interface for both multi-core CPUs and GPUs
- ▶ www.tnl-project.org
- ▶ MIT licence
- ▶ affiliated project of Numfocus (www.numfocus.org)



TEMPLATE
NUMERICAL
LIBRARY



Installation of TNL

The library can be downloaded using git as follows:

```
1 git clone gitlab@mmg-gitlab.fjfi.cvut.cz:tnl/tnl-dev.git tnl-dev
```

TNL is header-only library and so installation is very fast:

```
1 cd tnl-dev  
2 ./install
```

This installs TNL into \${HOME}/.local/include.

Working with address spaces

- ▶ CPU has its own system memory
 - ▶ GPU has its own global memory
 - ▶ both are connected by slow PCI Express bus
 - ▶ programmer must carefully distinct the two address spaces

Memory allocation in TNL

Memory allocation is done by templated class Vector:

```
1  namespace TNL::Containers;
2  template< typename Real = double,
3            typename Device = TNL::Devices::Host,
4            typename Index = int >
5  class Vector { ... };
```

where

- ▶ Real is type of elements stored in the vector
- ▶ Device says where the vector will be allocated
 - ▶ TNL::Devices::Host for CPU and system memory
 - ▶ TNL::Devices::Cuda for GPU and global memory
- ▶ Index is type for indexing the elements in the vector

Memory allocation in TNL

```
1 #include <iostream>
2 #include <list>
3 #include <vector>
4 #include <TNL/Containers/Vector.h>
5
6 using namespace TNL;
7 using namespace TNL::Containers;
8
9 int main( int argc, char* argv[] )
10 {
11     Vector< int > host_vector( 10 );           // vector with 10 elements on CPU
12     Vector< int, Devices::Cuda > device_vector; // empty vector on GPU
13
14     host_vector = 3;                          // set all elements to 3
15     device_vector = host_vector;             // copy the host vector on GPU
16
17     std::cout << "host_vector = " << host_vector << std::endl;
18     std::cout << "device_vector = " << device_vector << std::endl;
19     std::cout << std::endl;
```

Memory allocation in TNL

```
1     std::list< int > list { 1, 2, 3, 4, 5 };
2     std::vector< int > vector { 6, 7, 8, 9, 10 };
3
4     Vector< int, Devices::Cuda > device_vector_list( list );
5     Vector< int, Devices::Cuda > device_vector_vector( vector );
6     Vector< int, Devices::Cuda > device_vector_init_list{ 11, 12, 13, 14, 15 };
7
8     std::cout << "device_vector_list = " << device_vector_list << std::endl;
9     std::cout << "device_vector_vector = " << device_vector_vector << std::endl;
10    std::cout << "device_vector_init_list = " << device_vector_init_list << std::endl;
11 }
```

The result looks as follows:

```
1 host_vector = [ 3, 3, 3, 3, 3, 3, 3, 3, 3, 3 ]
2 device_vector = [ 3, 3, 3, 3, 3, 3, 3, 3, 3, 3 ]
3
4 device_vector_list = [ 1, 2, 3, 4, 5 ]
5 device_vector_vector = [ 6, 7, 8, 9, 10 ]
6 device_vector_init_list = [ 11, 12, 13, 14, 15 ]
```

Expression templates for vectors

The following vector expression

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

can be evaluated with Blas/Cublas as follows:

```
1 cublasHandle_t c_h;
2 cublasSaxpy(c_h,N,1.0,a,1,x,1); // -> x = a (assume x = [0,...,0] at the beginning)
3 cublasSaxpy(c_h,N,2.0,b,1,x,1); // -> x = x + 2 * b
4 cublasSaxpy(c_h,N,3.0,c,1,x,1); // -> x = x + 3 * c
```

And using the **expression templates (ET)** in TNL as follows:

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

The later is **more efficient**.

Expression templates for vectors

- ▶ ET is a proxy object for vector expression
- ▶ it allows to do **lazy evaluation of the expression**
- ▶ it can evaluate i -th element of the expression on-the-fly, i.e.

```
1     ET[ i ] = a[ i ] + 2 * b[ i ] + 3 * c[ i ]
```

- ▶ on CPU, the assignment would be evaluated as follows

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

- ▶ there is only **one write** to \vec{x} instead of **3 writes and 2 reads** in case of Blas

Expression templates for vectors

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

Size	CPU			GPU		
	BLAS	TNL		cubLAS	TNL	
	BW	BW	Speed-up	BW	BW	Speed-up
100k	19.3	41.5	2.2	194.7	236.5	1.21
200k	19.7	41.7	2.1	228.3	277.6	1.21
400k	17.3	35.9	2.1	218.3	330.9	1.51
800k	11.7	19.3	1.6	233.3	370.6	1.58
1.6M	10.4	17.0	1.6	249.6	403.4	1.61
3.2M	10.2	17.3	1.7	266.6	444.8	1.66
6.4M	10.2	17.3	1.7	276.6	471.3	1.70

BW = effective memory bandwidth in GB/s,

tested on GPU Nvidia P100 (16 GB HBM2 @ 732 GB/s, 3584 CUDA cores)
and Intel Core i7-5820K (3.3GHz, 16MB cache).

Vector operations

We setup the following vectors:

```
1  using namespace std;
2  using Vector = TNL::Containers::Vector< float, TNL::Devices::Cuda >;
3  Vector a{ 8, 4, 2, 0, -2, -4, -8 };
4  Vector b{ 0, 2, 4, 6, 8, 10, 12 };
```

We may search for maximum and minimum:

```
1  cout << "min( a ) = " << min( a ) << endl
2  << "max( a ) = " << max( a ) << endl;
```

The result looks as follows:

```
1  min( a ) = -8;
2  max( a ) = 8
```

Vector operations

We may combine vector operations with ET:

```
1 cout << "abs( a ) = " << abs( a ) << endl
2     << "min( abs( a ) ) = " << min( abs( a ) ) << endl
3     << "max( abs( a ) ) = " << max( abs( a ) ) << endl;
```

The result looks as follows:

```
1 abs( a ) = [ 8, 4, 2, 0, 2, 4, 8 ]
2 min( abs( a ) ) = 0
3 max( abs( a ) ) = 8
```

Vector operations

We may compute minimum and maximum componentwise:

```
1 cout << "min( a, b ) = " << min( a, b ) << endl
2     << "max( a, b ) = " << max( a, b ) << endl
3     << "min( max( a, b ) ) = " << min( max( a, b ) ) << endl
4     << "max( abs( a ), b ) = " << max( abs( a ), b ) << endl;
```

The result looks as follows:

```
1 min( a, b ) = [ 0, 2, 2, 0, -2, -4, -8 ]
2 max( a, b ) = [ 8, 4, 4, 6, 8, 10, 12 ]
3 min( max( a, b ) ) = -8
4 max( abs( a ), b ) = [ 8, 4, 4, 6, 8, 10, 12 ]
```

Vector operations

We may locate minimal and maximal elements:

```
1 auto arg_min_a = argMin( a );      // -> std::pair< float, int >
2 auto arg_max_a = argMax( a + b ); // -> std::pair< float, int >
3 cout << "min( a ) = " << arg_min_a.first << " at " << arg_min_a.second << endl
4     << "max( a + b ) = " << arg_max_a.first << " at " << arg_max_a.second << endl;
```

The result looks as follows:

```
1 min( a ) = -8 at 6
2 max( a + b ) = 8 at 0
```

Vector operations

We may perform componentwise operations:

```
1 cout << "a + b = " << a + b << endl  
2     << "a - b = " << a - b << endl  
3     << "a * b = " << a * b << endl  
4     << "a / b = " << a / b << endl;
```

The result looks as follows:

```
1 a + b = [ 8, 6, 6, 6, 6, 6, 4 ]  
2 a - b = [ 8, 2, -2, -6, -10, -14, -20 ]  
3 a * b = [ 0, 8, 8, 0, -16, -40, -96 ]  
4 a / b = [ inf, 2, 0.5, 0, -0.25, -0.4, -0.666667 ]
```

Vector operations

We may compute norms and scalar products:

```
1 cout << "l2Norm( a - b ) = " << l2Norm( a - b ) << endl
2     << "l2Norm( a + 3 * sin( b ) ) = " << l2Norm( a + 3 * sin( b ) ) << endl
3     << "Scalar product: ( a, b ) = " << ( a, b ) << endl
4     << "Scalar product: ( a + 3, abs( b ) / 2 ) = " << ( a + 3, abs( b ) / 2 )
5     << endl;
```

The result looks as follows:

```
1 l2Norm( a - b ) = 28.3549
2 l2Norm( a + 3 * sin( b ) ) = 15.3312
3 Scalar product: ( a, b ) = -136
4 Scalar product: ( a + 3, abs( b ) / 2 ) = -5
```

Using lambda functions

More complex operations can be done with **lambda functions**.

```
1 using Vector = TNL::Containers::Vector< float, TNL::Devices::Cuda >;
2 Vector a( 6, 1.0 );
3 cout << "a = " << a << endl;
4
5 auto f = [] __cuda_callable__ ( int idx, float& value ) { // This all is
6     value += 0.5 * idx;                                // turned into
7 };
8 a.forAllElements( f );                                // a CUDA kernel
9
10 cout << "a = " << a << endl;                      // by C++ compiler.
```

The result looks as:

```
1 a = [ 1, 1, 1, 1, 1, 1 ]
2 a = [ 1, 1.5, 2, 2.5, 3, 3.5 ]
```

Computations on GPU

The lambda functions can **capture** variables from the surrounding code:

```
1  using Vector = TNL::Containers::Vector< float, TNL::Devices::Cuda >;
2  Vector a( 6, 1.0 );
3  cout << "a = " << a << endl;
4
5  const float h = 0.5; // capture this variable h
6  auto f = [=] __cuda_callable__ ( int idx, float& value ) {
7      value += h * idx; // use it here inside the lambda function
8  };
9  a.forAllElements( f );
10 cout << "a = " << a << endl;
```

We do not have to copy all necessary variables on GPU explicitly.

Computations on GPU

- ▶ if the lambda function is executed on GPU, the captured **variable must be copied on the GPU**
- ▶ this is why all variables **must be captured as a copy** not as a reference
 - ▶ we must use the statement [=] in the definition of the lambda function
 - ▶ we cannot use [&]
- ▶ with the lambda functions, the captured variables are transferred on the GPU **automatically**
- ▶ if we write CUDA kernels, we have to transfer all variables explicitly = very annoying

Vector view

- ▶ what if we need to work with **two vectors**?
- ▶ vector **cannot be captured as a copy**:
 - ▶ it would create a new copy of the vector ⇒ **additional memory** is needed
 - ▶ we would not be able **to modify the original vector**, only the copy
 - ▶ the copy is temporal and it is lost when the lambda function finishes
- ▶ we have to use a **vector view** instead of vector
 - ▶ vector view shares data with another vector
 - ▶ copy of a vector view is not a deep copy ⇒ elements of the original vector are not duplicated

Vector view - example

```
1  using Vector = TNL::Containers::Vector< float, TNL::Devices::Cuda >;
2  Vector a( 6, 1.0 ), b( 6, 3.0 );
3  cout << "a = " << a << endl;
4
5  const float h = 0.5;
6  auto b_view = b.getView();           // this view can be captured by making copy
7  auto f = [=] __cuda_callable__ ( int idx, float& value ) {
8      value = h * idx + b_view[ idx ]; // we can use the view even on the GPU
9  };
10 a.forAllElements( f );
11 cout << "a = " << a << endl;
```

The result looks as follows:

```
1  a = [ 1, 1, 1, 1, 1, 1 ]
2  a = [ 3, 3.5, 4, 4.5, 5, 5.5 ]
```

Vector view - example

Vector view can be used for encapsulating data not allocated by TNL:

```
1  using VectorView = TNL::Containers::VectorView< float, TNL::Devices::Host >;
2  float* data = new float[ 6 ];
3  VectorView a( data, 6 );
4  a = 1.0;
5  cout << "a = " << a << endl;
6
7  auto f = [] __cuda_callable__ ( int idx, float& value ) {
8      value += 0.5 * idx;
9  };
10 a.forAllElements( f );
11 cout << "a = " << a << endl;
12
13 delete[] data;
```

Reduction

Reduction is an operation that takes 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 = 100 lines of code

Flexible parallel reduction

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 auto fetch = [=] __cuda_callable__ (int i)->float { return a[i]*b[i]; };
2 auto reduction = [] __cuda_callable__ (float x, float y) -> float { return x+y; };
3
4 float result( 0.0 );
5 for( int i = 0; i < size; i++ )
6     reduction = reduction( result, fetch( i ) );
```

Flexible parallel reduction

In TNL, the for-loop is replaced with call of `reduce` function:

```
1 auto a_view = a.getView();
2 auto b_view = b.getView();
3 auto fetch = [=] __cuda_callable__ ( int i ) -> float {
4     return a_view[ i ] * b_view[ i ]; };
5 auto reduction = [] __cuda_callable__ ( float x, float y ) -> float { return x + y; };
6
7 result = TNL::Algorithms::reduce< Device >( 0, a.getSize(), fetch, reduction, 0.0 );
```

The last parameter is the identity element for given reduction operation.

Flexible parallel reduction

Comparison of two vectors can be evaluated as follows:

```
1 auto a_view = a.getView();
2 auto b_view = b.getView();
3 auto fetch = [=] __cuda_callable__ (int i)->bool {
4     return a_view[ i ] == b_view[ i ]; };
5 auto reduction = [] __cuda_callable__ (float x, float y)->float { return x && y; };
6
7 result = TNL::Algorithms::reduce< Device >( 0, a.getSize(), fetch, reduction, true );
```

Flexible parallel reduction

Largest element in absolute value can be evaluated as follows:

```
1 auto a_view = a.getView();
2 auto fetch = [=] __cuda_callable__ (int i)->float { return abs( a_view[ i ] ); };
3 auto reduction = [] __cuda_callable__ (float x, float y)->float {
4     return max( x, y ); };
5
6 result = TNL::Algorithms::reduce< Device >( 0, a.getSize(), fetch, reduction,
7             std::numeric_limits< float >::lowest() );
```

Or more compact way:

```
1 auto a_view = a.getView();
2 result = TNL::Algorithms::reduce< Device >( 0, a.getSize(),
3             [=] (int i)->float { return abs( a_view[ i ] ); },
4             TNL::Max() );
```

Flexible parallel reduction

We can merge two operations together:

- ▶ update/addition of a vector
- ▶ computation of the norm of the update

It appears often in Runge-Kutta solvers.

```
1 auto a_view = a.getView();
2 auto b_view = b.getView();
3 result = TNL::Algorithms::reduce< Device >( 0, a.getSize(),
4                                     [=] __cuda_callable__ (int i)->float mutable {
5                                         // mutable because we change a_view
6                                         float update = 0.5 * ( a_view[ i ] + b_view[ i ] );
7                                         a_view[ i ] += update;
8                                         return abs( update ); },
9                                     TNL::Plus() );
```

Solving the heat equation

At the end, we show FDM solver for the heat equation given as:

$$\begin{aligned}\frac{\partial u(\vec{x}, t)}{\partial t} - \Delta u(\vec{x}, t) &= 0 \text{ on } \Omega \times [0, T], \\ u(\vec{x}, 0) &= u_{ini}(\vec{x}) \text{ on } \Omega, \\ u(\vec{x}, t) &= 0 \text{ on } \partial\Omega \times [0, T].\end{aligned}$$

We approximate it by the finite difference method:

$$\begin{aligned}u_{ij}^{k+1} &= u_{ij}^k + \frac{\tau}{h^2} (u_{i-1,j} + u_{i,j-1} + u_{i+1,j} + u_{i,j+1} - 4u_{ij}) \text{ on } \Omega_h, \\ u_{ij}^{k+1} &= 0 \text{ on } \partial\Omega_h.\end{aligned}$$

Solving the heat equation

```
1 #include <TNL/Containers/Vector.h>
2 #include <TNL/Algorithms/reduce.h>
3 #include <TNL/Algorithms/ParallelFor.h>
4
5 using Device = TNL::Devices::Host;
6
7 int main( int argc, char* argv[] )
8 {
9     using Vector = TNL::Containers::Vector< float, Device >;
10
11     // Parameters of the discretization
12     int size = 11;
13     float h = 1.0 / ( size - 1 );
14     float tau = 0.1 * h * h;
15     float final_T = 1.0;
16     float h_inv_sqr = 1.0 / ( h * h );
```

Solving the heat equation

```
18     /////
19     // Allocation of the grid functions
20     Vector u( size * size, 0.0 ), aux( size * size, 0.0 );
21     auto u_view = u.getView();
22     TNL::Algorithms::ParallelFor2D< Device >::exec( 0, 0, size, size,
23         [=] __cuda_callable__ ( int i, int j ) mutable {
24             float x = i * h - 0.5;
25             float y = j * h - 0.5;
26             int idx = j * size + i;
27             if( x * x + y * y < 0.25 )
28                 u_view[ idx ] = 1.0; } );
```

Solving the heat equation

```
29     //////
30     // Time loop
31     float t = 0.0;
32     float residue = 10.0;
33     while( t < final_T && residue > 1.0e-6 )
34     {
35         //////
36         // Update with FDM approximation of the Laplace operator
37         auto u_view = u.getView();
38         auto aux_view = aux.getView();
```

Solving the heat equation

```
39     residue = sqrt( TNL::Algorithms::reduce< Device >( 0, size * size,
40                     [=] __cuda_callable__ ( int idx ) mutable -> float {
41                         int i = idx % size;
42                         int j = idx / size;
43                         if( i == 0 || j == 0 || i == size - 1 || j == size - 1 )
44                             return 0.0;
45                         float update = h_inv_sqr * (
46                             u_view[ idx - size ] + u_view[ idx - 1 ] +
47                             u_view[ idx + size ] + u_view[ idx + 1 ]
48                             - 4.0 * u_view[ idx ] );
49                         aux_view[ idx ] = u_view[ idx ] + tau * update;
50                         return update * update;
51                     },
52                     TNL::Plus() ) );
```

Solving the heat equation

```
53         t += tau;
54         u.swap( aux );
55         std::cout << "t = " << t << " residue = " << residue << std::endl;
56     }
57     return EXIT_SUCCESS;
58 }
```

Summary

- ▶ TNL allows simple and efficient way how to develop code for GPUs
- ▶ there is no performance drop
- ▶ the code can be written and **DEBUGGED** on CPU
- ▶ with a little bit of luck it runs on GPU as well