

Einleitung

Um deal.II Solver Klassen mit CUDA zu verwenden, werden Vektoren, die speziell für CUDA implementiert sind, benötigt. Die CUDA Vektoren müssen die Vektor-Schnittstelle von deal.II erfüllen, also alle benötigten Methoden implementieren.

Die deal.II Vektorschnittstelle sieht wie folgt aus:

```
1 class VECTOR {
2     public:
3
4     void reinit (const VECTOR&, bool leave_elements_uninitialized = false);
5
6     double operator * (const VECTOR &v) const;
7
8     void add (const VECTOR &v);
9
10    void add (const double a, const VECTOR &v);
11
12    void sadd (const double a, const double b, const VECTOR &v);
13
14    void equ (const double a, const VECTOR &v);
15
16    VECTOR & operator += (const double a);
17
18    double l2_norm () const;
19};
```

Des Weiteren müssen die Speicherformatklassen die Matrix-Schnittstelle aus deal.II erfüllen:

```
1 class Matrix {
2     public:
3     void vmult (VECTOR &dst, const VECTOR &src) const;
4
5     void Tvmult (VECTOR &dst, const VECTOR &src) const;
6
7 };
8
```

Ziel

Ziel von Step-6 war die Implementierung einer Vektorklasse, die die Vektoroperationen mit Hilfe von CUDA auf der GPU ausführt.

Dabei sollten nur die von den deal.II Solver-Klassen benötigten Methoden implementiert werden.

Die Hauptschwierigkeit bestand in der korrekten Allokation auf dem Grafikkartenspeicher.

Implementierung

Bei der Implementierung werden die Vektorwerte, seien dies Standard- oder deal.II Vektoren, die auf dem Host liegen, auf die Grafikkarte gespiegelt, indem Zeiger auf diese Daten der GPU übergeben werden und damit auch Speicher auf der GPU allokiert wird.

Des Weiteren werden Kernel für die elementaren Vektoroperationen implementiert, die dies auf der GPU ausführen. Die Ergebnisse werden dann auf den Host zurück kopiert, um die Weiterbearbeitung bzw. Ausgabe zu ermöglichen.

Beispiel Kernel: Addition

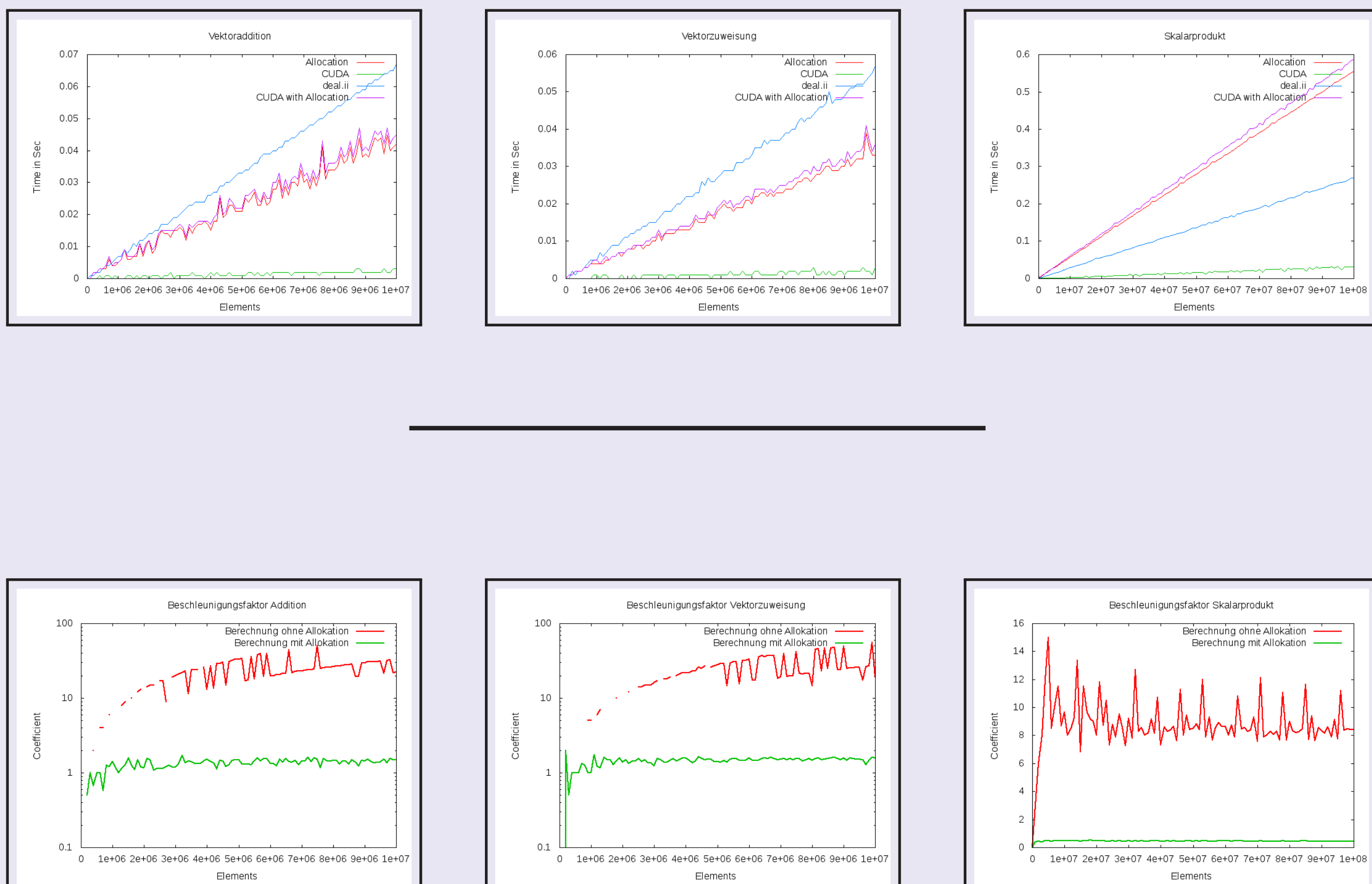
```
1 template <typename T>
2 __global__ void VecAdd(T* A, T* B, int N)
3 {
4     int threadId = threadIdx.x + blockDim.x * blockIdx.x;
5     int gridSize = blockDim.x * blockIdx.x;
6
7     for(int i = threadId; i < N; i += gridSize){
8         A[i] += B[i];
9     }
10 }
11
12 // gekapselter Kernelaufruf
13 void call_kernel_VecAdd(int _NUM_BLOCKS, int BLOCK_SIZE, double* A, double* B,
14                        int N){
15     VecAdd<<_NUM_BLOCKS, BLOCK_SIZE>>(A, B, N);
16     cudaThreadSynchronize();
17 }
```

Aufruf aus C++ Programm :

```
1 // Treiber fuer den GPU-Teil
2 #include "cuda_driver.h"
3 #include <cutil_inline.h>
4 #include "cuda_driver.cu.hh"
5 ...
6
7 // l&ouml;ser aus deal.II mit GPU-Vektor initialisieren
8 SolverCG<CUDAVectorView<double> > cg (solver_control);
9
10 // ... um damit die GPU-seitigen Vektoren zu initialisieren
11 CUDAVectorView<double> solview (solution);
12 CUDAVectorView<double> rhsview (system_rhs);
13
14 // ... um unsere GPU-seitige Matrix damit zu initialisieren
15 CudaEIIItMatrixView<double> sm_cuda_view (system_matrix);
16
17 // ... PC initialisieren (geht genauso)
18
19 // Problem l&ouml;sen
20 cg.solve (sm_cuda_view, solview, rhsview, preconditioner_cuda);
21
22 solview.copyToHost();
```

Ergebnisse

Im Vergleich zur CPU Lösung mit deal.II, können folgende Resultate für elementare Vektoroperationen betrachtet werden:



Beobachtungen

Da die Klasse CUDAVectorView als Grundlage für folgenden Klassen aus step-7 und step-8 dient, können die Resultate bei den Klassen CudaCSR-MatrixView, CudaEIIItMatrixView und CudaEIIIRMatrixView betrachtet werden.

- Anhand der Abbildungen erkennt man, dass die Rechenzeit linear zur Problemgröße wächst.
- Im Falle der Vektoraddition und der Zuweisung wird der Vorteil von CUDA recht deutlich im Vergleich zur CPU-Lösung. So erhält man doppelt so schnell eine Lösung, betrachtet man die Allokation jedoch nicht so hat man einen Performance-Vorteil von Faktor 6.
- Im Falle der Berechnung des Skalarprodukts wird deutlich, dass die Rechenzeit mit CUDAVectorView Vektoren viel niedriger ist als mit deal.II, beachtet man aber die Allokationszeit, so ist CUDAVectorView weniger effizient im Vergleich zur deal.II. Da die Vektoren nur einmal allokiert werden, kann man davon ausgehen, dass bei Berechnungen, die immer wieder die selben Vektoren benutzen, mit CUDA die schnellere Rechenzeit erhält.

Literatur

- Yousef Saad - Iterative methods for sparse linear systems (2nd edition).
- Numerische Mathematik, interaktiv Uni-Wien (<http://www.dorn.org/uni/sls/index.html>).
- deal.II Dokumentation (<http://www.dealii.org/6.2.1/doxygen/tutorial/index.html>)