

1 Einführung und Allgemeines

In diesem Praktikum sollten verschiedene parallele Zugriffsmuster auf Grafikkarten gemessen und verglichen werden. Dabei verwendet wurde CUDA [1], eine von Nvidia entwickelte C++-Erweiterung, die Grafikkarten dieser Firma für parallelisierbare Berechnungen abseits von klassischem Rendering nutzbar macht. Konkret durchgeführt wurden die Messungen auf einer GTX 1070. Einige relevante Spezifikationen dieser GPU sind:

- Cores: 1920, wovon ein Warp jeweils 32 Threads umfasst [2]
- Maximale Taktung laut „nvidia-smi“: 1708 MHz [2]
- Speicherkonfiguration: 8 GB GDDR5 RAM [2]
- Maximaler Durchsatz (Herstellerangabe): 256 GB/s [2]
- PCI: PCIe 3.0 fähig, in der Messung wurde allerdings nur PCIe 2.0x16 verwendet. Die maximale Kopierrate von PCIe 2.0x16 entspricht 8.0GB/s [3]

2 Grundsätzliche Beobachtungen

Eine Grundlage der Programmierung mit der CUDA ist das Management der Daten, welche auf der CPU (genannt „Host“) oder der GPU (genannt „Client“) liegen, da CPU-Kerne nur bedingt auf den GPU-RAM zugreifen können und umgekehrt. Mithin ist die mit Standard-Kopieroperatoren erzielbare Kopierate interessant, die gleichzeitig auch einen guten Vergleichsmaßstab für die Effizienz der verschiedenen Kernel darstellt. Abbildung 1 zeigt die für `cudaMemcpy()` gemessenen Kopieraten und die vom Hersteller angegebenen Maximalwerte sowie die Spezifikation für den verwendeten PCIe.

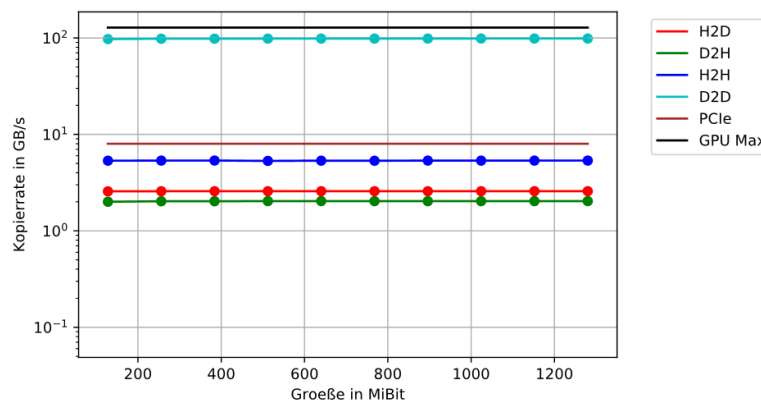


Abbildung 1: Kopierate von `cudaMemcpy()`

Bei der tatsächlichen Durchführungen der Messungen muss beachtet werden, dass jeweils der erste Kernel Aufruf eines CUDA Programms eine wesentlich längere Zeit be-

nötigt als darauf folgende Aufrufe des selben Kernels. Abbildung 2 zeigt eine Beispielmessung für einen Empty Kernel, der zwar auf der GPU ausgeführt wird, dort allerdings keine Arbeit verrichtet. Um dieser Verzerrung entgegenzuwirken, muss vor jeder Laufzeitmessung für einen Kernel dieser einmal als „Warmup“ aufgerufen, aber dann aus der eigentlichen Messung herausgerechnet werden.

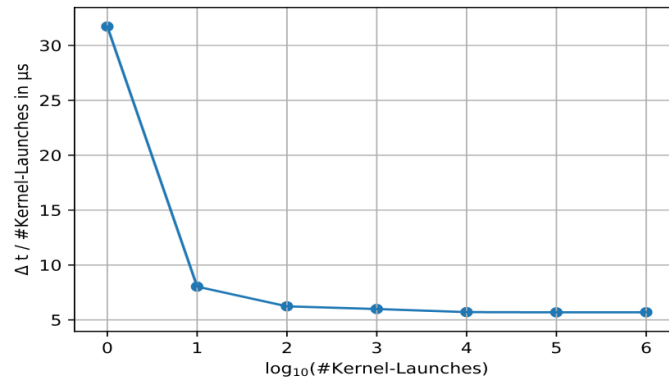


Abbildung 2: Startup Cost eines Empty Kerneles

3 Copy Kernel

Der „Copy Kernel“, der jeweils ein Element eines Buffers in die entsprechende Position eines anderen Buffers kopiert, lässt sich durch folgenden CUDA-Quellcode beschreiben:

```
template<typename T>
__global__
void copyKernel(T* out, T* in) {
    unsigned id = threadIdx.x + blockIdx.x * blockDim.x;

    out[id] = in[id];
}
```

Obwohl die Funktion als Input nur Zeiger auf einen Quell- und einen Ziel-Buffer erhält, gibt es doch drei implizite Parameter. Diese treten auch bei allen folgenden Messungen auf:

1. Die Anzahl der Blöcke, in die die einzelnen Threads von der GPU gruppiert werden
2. Die Größe dieser Blöcke
3. Den Zugriffstyp T (zB. `char`, `int`)

Abbildung 3 und Abbildung 4 zeigen die Auswirkungen von Variationen über die drei Parameter auf die durch den Copy Kernel erzielte Kopierrate. Dabei zeigt sich, dass

sowohl für die Anzahl der Blöcke, für die Blockgröße als auch für `sizeof(T)` ein jeweils höherer Parameter zu einer höheren Kopierrate führt. Für hinreichend hohe Werte wird eine Sättigung nahe des Ergebnisses aus Abbildung 1 erzielt. Interessant ist außerdem, dass die Abbildungen sich für steigendes `sizeof(T)` je linearer verhalten, je geringer Blockgröße bzw. Blockanzahl sind.

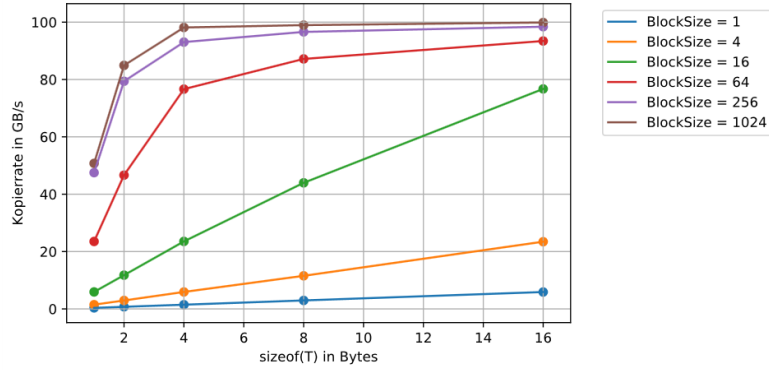


Abbildung 3: Copy Kernel: Variation über `BlockSize` und `sizeof(T)`, `numBlocks` = 16384

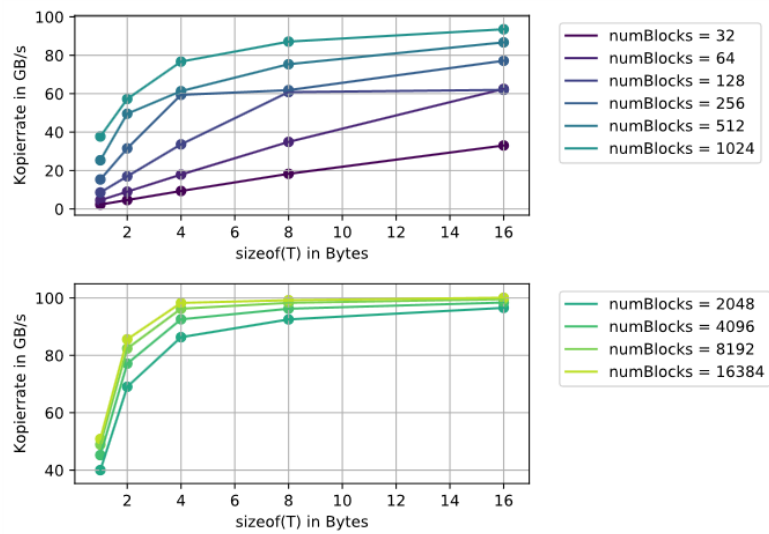


Abbildung 4: Copy Kernel: Variation über `numBlocks` und `sizeof(T)`, `BlockSize` = 1024

4 Strided Access

Der Strided Access ist wie folgt definiert:

```

template<typename T>
__global__
void copyKernel(T* out, T* in, int stride) {
    unsigned id = threadIdx.x + blockIdx.x * blockDim.x;

    out[id*stride] = in[id*stride];
}

```

Das heißt, der Zugriff erfolgt nicht mehr auf jedes Element konsekutiv hintereinander, sondern nur auf jedes N -te Element. Zusätzlicher Parameter ist dann der **stride**, der festlegt, wie viele Buffer-Einträge übersprungen werden. Bei elementarem Zugriff, wie er zum Beispiel in einer Turingmaschine definiert ist, sollte sich keine Änderung der Kopiererrate ergeben (sofern die nicht bearbeiteten Elemente nicht mitgerechnet werden). Bei der Messung zeigt sich allerdings, dass deutliche Cashing-Effekte auftreten. Bei Messungen zeigt sich ein nahezu exponentieller Abfall der Kopiererrate mit steigendem **stride**, für den die beiden Messungen Abbildung 5 und Abbildung 6 beispielhaft stehen. Bei $\text{sizeof}(T) \cdot \text{stride} = 32$ und $\text{sizeof}(T) \cdot \text{stride} = 64$ ergeben sich jeweils zwei fast gleiche Werte, und ab $\text{sizeof}(T) \cdot \text{stride} = 2048$ fällt die Kopiererrate wenn überhaupt nur noch vernachlässigbar ab (bei diesen beiden Messungen wie auch bei allen anderen durchgeführten).

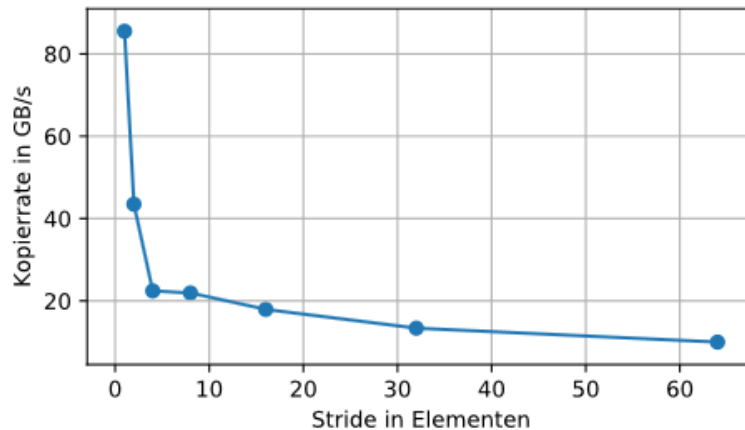


Abbildung 5: Strided Access: Variation über **stride**, BlockSize = 256, numBlocks = 4096, T: int2 (sizeof(int2) = 8)

5 Offset Access

Der Offset Access ist ein etwas komplexerer Kernel. Sein Quellcode lautet wie folgt:

```

template<typename T>
__global__

```

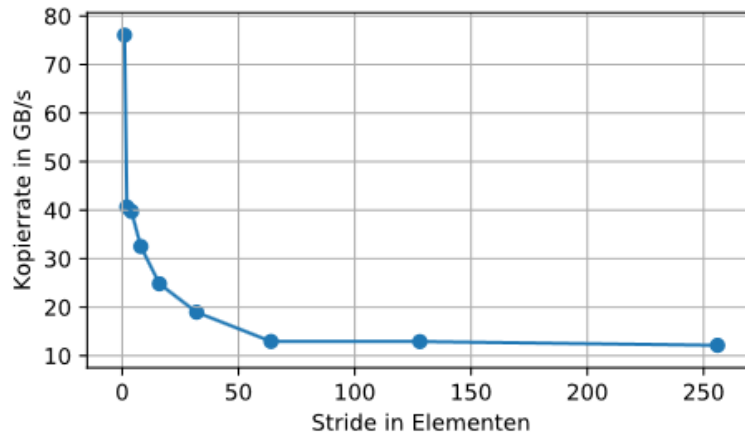


Abbildung 6: Strided Access: Variation über `stride`, `BlockSize = 128`, `numBlocks = 2048`, `T: int4` (`sizeof(int4) = 16`)

```

void oftKernel (T* out ,
                T* in ,
                const unsigned int sd_size ,
                const unsigned int block_size ,
                const unsigned int I ,
                const unsigned int L)
{
    const unsigned int sd_id = static_cast<int> (threadIdx.x / L);
    const unsigned int id = threadIdx.x - sd_id * L;
    const unsigned int sd_start = blockIdx.x *
                                blockDim.x * I + sd_id * L * I;

    for (unsigned int i = 0; i < I; i++)
    {
        const unsigned int el_id = sd_start + i * L + id;
        ((T*) out)[el_id] = ((T*) in)[el_id];
    }
}

```

Abbildung 7 zeigt das Zugriffsmuster, das aus diesem Kernel resultiert. Ein „oftKernel“ greift nicht konsekutiv mit jedem Thread auf ein Element der Buffer zu. Er teilt die Buffer zwar in gleich große Abschnitte auf, die jeweils einem Block zugeordnet werden, unterteilt die Blöcke allerdings erneut in Subdomains. Diese Subdomains bestehen dann aus einem gewissen Anteil Threads, die jeweils eine gewisse Anzahl an Iterationen durchführen müssen, um alle nötigen Kopiervorgänge abzuarbeiten. Neben den bekannten Parametern `T` und `BlockSize` ergeben sich dann noch zusätzlich die Anzahl der Iterationen pro Thread `I` und die Anzahl der Threads pro Subdomain `L`. Die Anzahl der Blöcke ergibt sich dann

aus diesen Parametern und der Größe des jeweils einer Subdomain zugeordneten Anteils der Buffer `sd_size`.

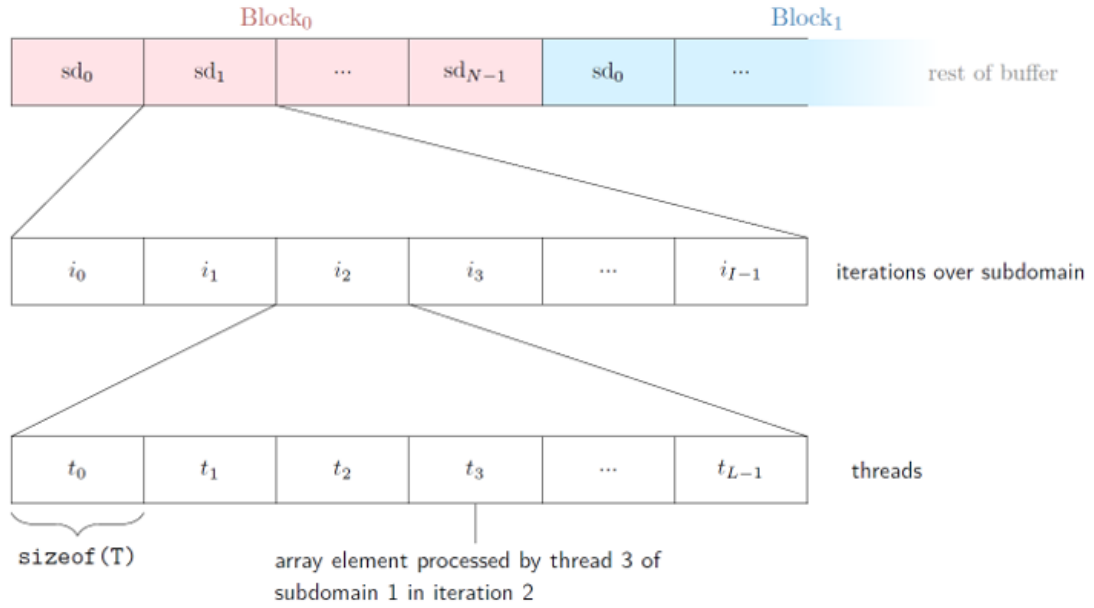


Abbildung 7: Zugriffsmuster des Offset Access

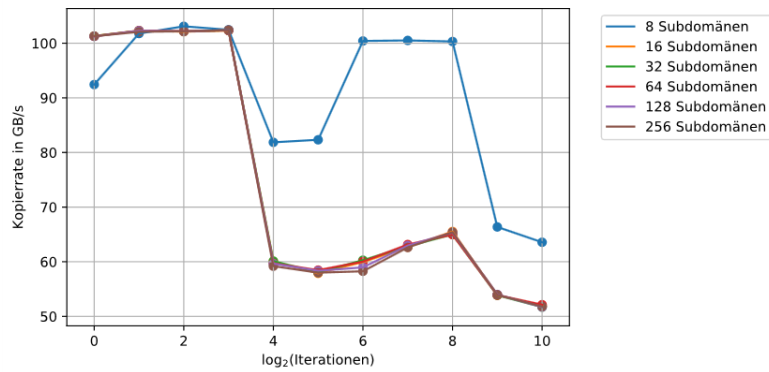


Abbildung 8: Offset Access: Variation über Subdomänen pro Block und $I, L = 4$, `sizeof(T) = 16`

Die Abbildungen 8, 9 und 10 zeigen verschiedene Variationen über die Parameter des Kernels. Es zeigt sich, dass die Zusammenhänge im Detail komplex und ohne eine tiefere Analyse der verschiedenen Caches der verwendeten GPU (isb. Gesamt- und Zugriffsgröße) wohl nicht vollständig zu verstehen sind. Dennoch lässt sich festhalten, dass eine Sättigung vergleichbar mit vorherigen Kopieroperationen erreicht werden kann, wenn die Parameter richtig gewählt werden. Insbesondere können auch Konfigurationen eine

Sättigung erreichen, bei denen ein Warp, d.h. die jeweils 32 immer gleich ausgeführten Threads, nicht auf nacheinander liegende Elemente zugreift. Siehe hierfür insbesondere Abbildung 10. Außerdem zeigt sich ein deutlicher Zusammenhang mit der L2-Cache Line-Size von 128 Byte – ist die Zugriffsgröße einer Iteration 128 Byte oder ein ganzzahliger Teiler davon, ist die Kopierrate wesentlich höher.

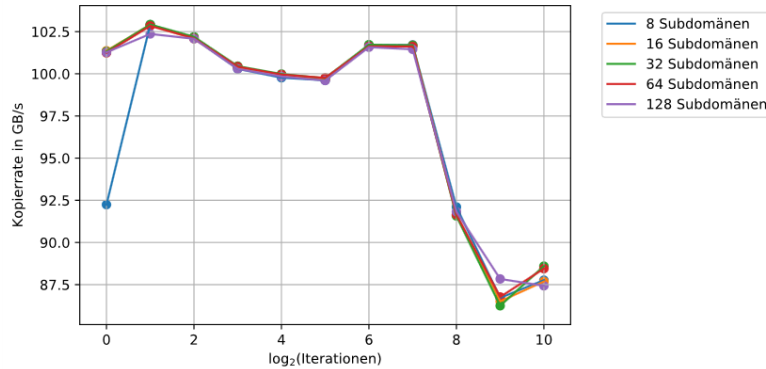


Abbildung 9: Offset Access: Variation über Subdomänen pro Block und I, $L = 8$, $\text{sizeof}(T) = 16$

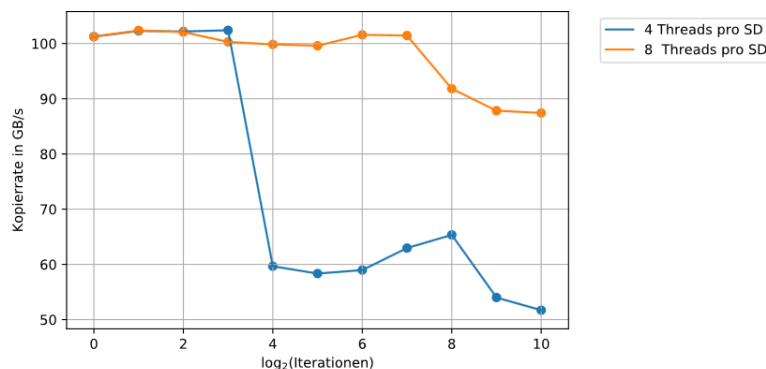


Abbildung 10: Offset Access: Variation über L und I, 128 Subdomänen pro Block, $\text{sizeof}(T) = 16$

6 Allokation: Standard und UnifiedMemory

Neben der Standard-Allokation über die C++ Funktion `malloc()` kann CUDA auch mit `cudaMallocHost()` einen sogenannten Pinned Memory Buffer auf dem Host-RAM allokiert. Dieser erlaubt schnelleres Kopieren und direkte Zugriffe der Kernel auf den Host-RAM, ohne vorher im Code händisch eine Kopieroperation auf den GPU-Speicher vornehmen zu müssen. Mit einem Copy Kernel ähnlich dem aus Kapitel 3 kann dann

erneut eine Kopierraten-Messung vorgenommen werden (Abbildung 11). Einen Unterschied bezüglich der Anzahl Blöcke scheint es nicht zu geben, die Kurven liegen perfekt aufeinander. Die Kopierrate scheint linear zu steigen, bis eine Sättigung erreicht wird, die etwas unter dem Niveau des Geschwindigkeit des verwendeten PCIe's liegt (ähnlich dem Ergebnis aus Abbildung 1. Diese Sättigung wird beim Kopieren vom Device zum Host bereits mit weniger Threads erreicht, und das erreichte Niveau ist ein wenig höher.

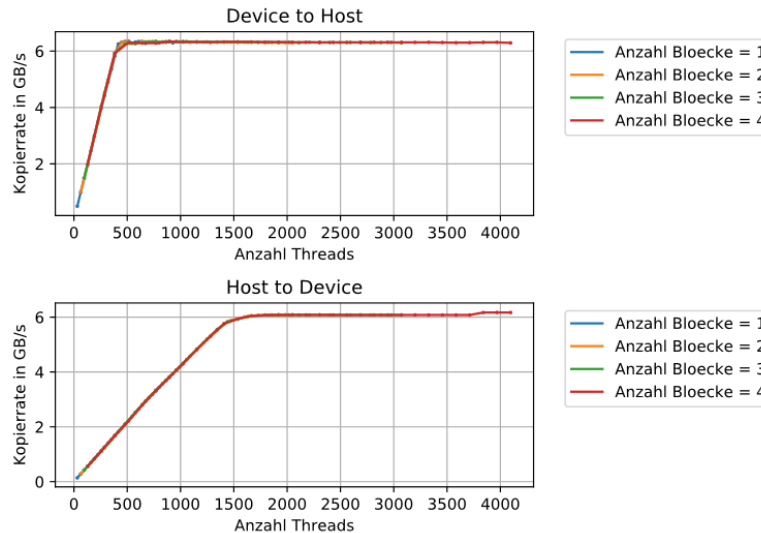


Abbildung 11: Copy Kernel Kopierrate auf Pinned Memory, `sizeof(T) = 4`

Zusätzlich können Buffer in CUDA auch mithilfe von `cudaMallocManaged()` als sogenannter Unified Memory allokiert werden. Dabei wird vom Programmierer überhaupt keine Angabe getroffen, ob der Buffer auf dem CPU- oder GPU-Ram liegt, sondern CUDA wickelt eventuell benötigte Kopiervorgänge automatisch ab. Abbildung 12 zeigt eine Kopierraten-Messung des Copy-Kernels auf Unified Memory. Dabei wurde durch Cs `memset()` und CUDAs `cudaMemset()` sichergestellt, dass die Input- und Output-Buffer vor der Messung jeweils vollständig auf den Host bzw. dem Device RAM kopiert wurden. Es zeigt sich ein ähnlicher Verlauf wie in der vorherigen Messung, allerdings ist die Form der Kurven viel ungleichmäßiger (möglicherweise ist der Messfehler größer), die Anzahl der Blöcke zeigt einen deutlicheren Unterschied und die Sättigung wird wesentlich später erreicht.

7 Kernel im Producer-Consumer-Verhältnis

Für diese Messung werden zwei identische Grafikkarten benötigt. Aufgrund der Komplexität des verwendeten Codes sei dessen Gestalt hier nur kurz umrissen: Zwei Kernel, die (in eine Richtung) miteinander kommunizieren müssen, d.h. der eine Kernel benötigt für die Erfüllung seiner Aufgabe die Ergebnisse des anderen aber nicht umgekehrt, stehen in einem Producer-Consumer-Verhältnis. Im vorliegenden Fall errechnet dabei der Producer

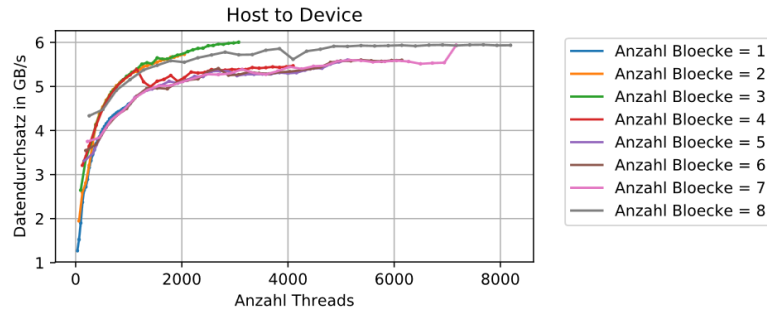


Abbildung 12: Copy Kernel Kopiertrate auf Unified Memory, T: int

den Durchschnitt eines Teils des Input-Buffers, während der Consumer diesen berechneten Durchschnitt ausliest und ihn vielfach in den entsprechenden Teil gleicher Größe des Output-Buffers dupliziert. Dabei wartet der Consumer jeweils nur auf einen Teil des Buffers, bevor er diesen dupliziert. Eine Grafikkarte führt dann jeweils einen Kernel aus. Eine Zeitmessung für das Producer-Consumer-System bei Ignorieren der Kommunikation findet sich in Abbildung 13: Die benötigte Zeit steigt nahezu linear mit den Ergebnissen pro Block, was auch zu erwarten ist, da mit diesem Parameter auch proportional die Menge der bearbeiteten Daten steigt.

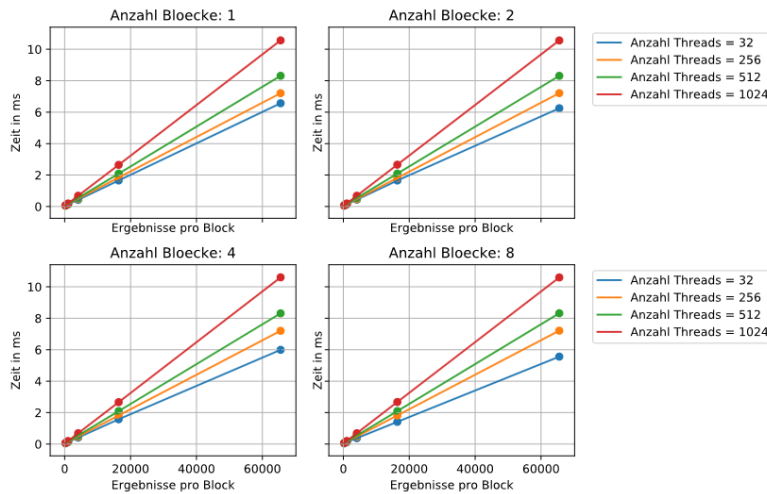


Abbildung 13: Laufzeit Producer-Consumer-Verhältnis ohne Kommunikation

Die Graphen mit aktivierter Kommunikation verlaufen ähnlich, allerdings auf einem wesentlich höheren Niveau. Demzufolge führt auch die Berechnung des Anteils der Kommunikation am Gesamtergebnis zu einem mehr oder weniger konstanten Verlauf der Kurven (Abbildung 14). Lediglich bei einer sehr geringen Anzahl von Ergebnissen pro Block zeigt sich eine gewisse Abweichung. Mit über 80%, teilweise sogar über 90% liegt der Kommunikationsanteil jedenfalls sehr hoch, und er scheint mit steigender Anzahl von

Blöcken auch leicht zu steigen. Umgekehrt handelt es sich bei dem Ausschlag der Kurve für 8 Blöcke und 1024 Threads beim letzten Wert vermutlich um einen Messfehler.

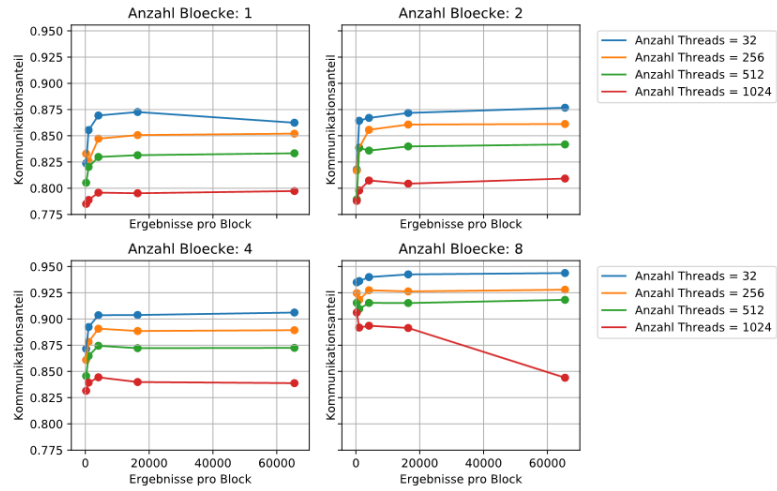


Abbildung 14: Anteil der durch Kommunikation verlorenen Zeit im Producer-Consumer-Verhältnis

Literatur

- [1] Nvidia Website: “CUDA Zone”, developer.nvidia.com/cuda-zone, aufgerufen am 6.2.2019
- [2] Nvidia Website: “GTX 1070”, nvidia.com/de-de/geforce/products/10series/geforce-gtx-1070, aufgerufen am 6.2.2019
- [3] Wikipedia: “PCI Express”, en.wikipedia.org/wiki/PCI_Express, aufgerufen am 6.2.2019