

---

# GPU-beschleunigtes Sortieren

---

Am Beispiel QuickSort auf einer  
Nvidia Grafikkarte

Von Michael Schmidt und Kevin Büchele

---

# Gliederung

1. Erklärung QuickSort
2. CUDA Hardwaremodell
3. Erklärung des Algorithmus`
4. Beispiel-Diagramme

---

# 1. QuickSort

- Finden eines Pivotelements
  - Zu sortierendes Array mit Pivot als „Seperator“ spalten
  - Rekursiv wird links vom Pivot und rechts vom Pivot nach selben Prinzip sortiert
  - Am Ende werden alle Teile zusammengefügt  
(trivialerweise ist danach kein Sortieren mehr notwendig)
-

---

## 2. CUDA Hardware

Definition:

- parallele Berechnungsarchitektur
- deutliche Steigerung der Rechenleistung

In anderen Worten:

- bietet Schnittstelle für Programmierern, um die Rechenleistung der GPUs mit einzubeziehen
  - insgesamt eine signifikante Steigerung der Rechenleistung zuerreichen.
-

---

## Weitere Eigenschaften (Nachteile)

- Meistens in C programmiert
  - Eher untypische Datentypen (9 Bit, 12 Bit)
    - Häufige Emulierung durch Software notwendig
  - Höhere Latenzzeiten durch Anbindung mittels PCI(e)
-

---

# Technische Eigenschaften der unterstützenden Grafikkarten

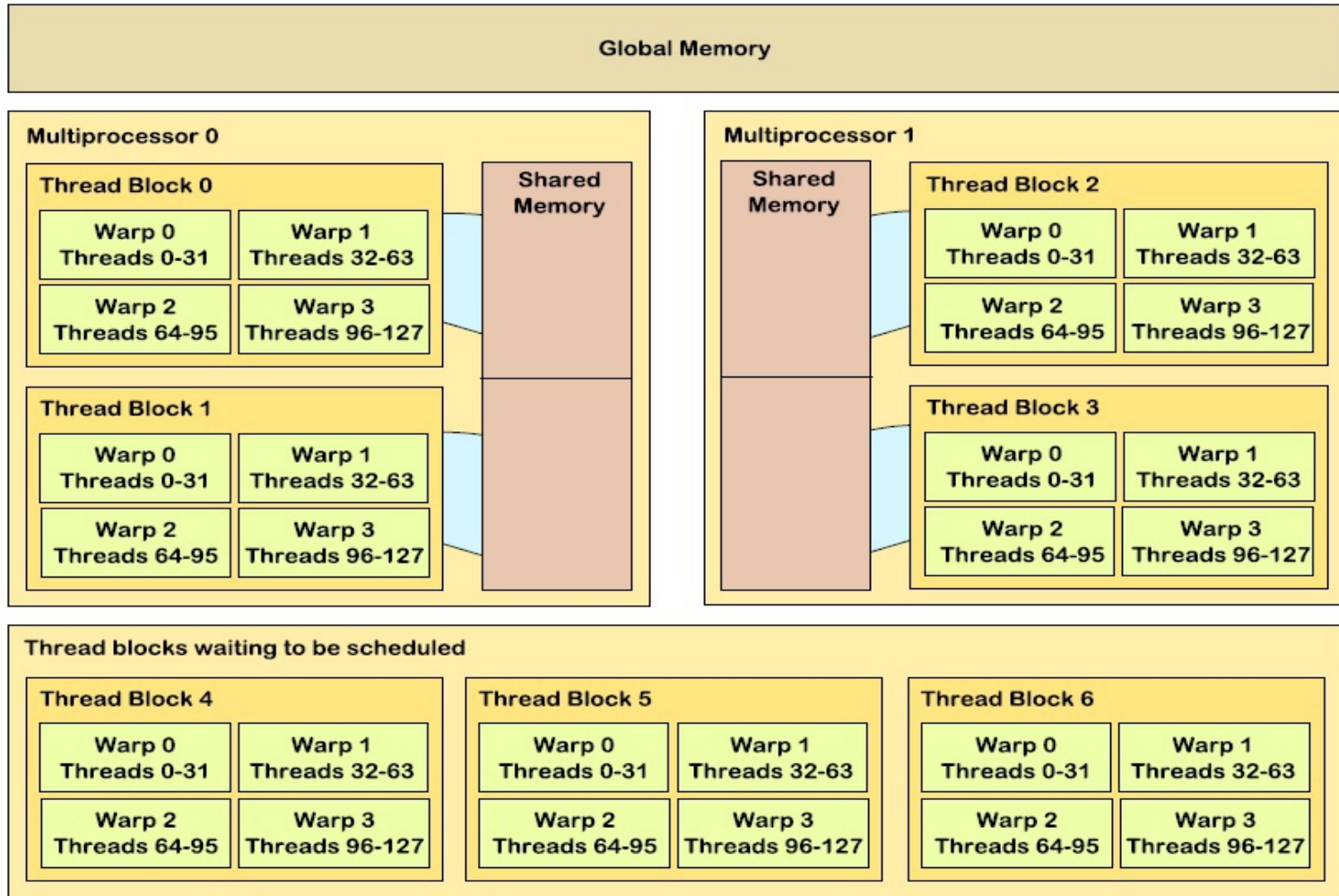
- 16 Multiprozessoren mit jeweils 8 Prozessoren
  - Jeder Multiprozessor unterstützt bis zu 768 Threads
  - 16KiB ( = 2048 Byte ) lokaler Speicher vorhanden mit maximal 8192 verfügbaren Registern
-

---

# Scheduling

- Threads zusammengefasst in Threadblocks
  - Abhängig von Anzahl der Register, benötigter Speicherplatz, können einzelne Threadblocks zum gleichen Multiprozessor zugeordnet werden
  - Warps fassen sequenziell/gleichzeitig arbeitende Threads zusammen
  - Es kann immer nur ein Warp pro Multiprozessor ausgeführt werden
-

# Grafische Darstellung der Hardware





# Leistung von CUDA

ISV	BESCHREIBUNG	VORTEILE DES GRAFIKPROZESSORS
<a href="#">MathWorks MATLAB®</a>	Datenparallele Berechnungen (MATLAB PCT, MDCS)	Deutliche Beschleunigung für mehr Produktivität für Studenten, Wissenschaftler und Techniker
<a href="#">Jedox Palo</a>	Excel-Erweiterung mit OLAP für Planung und Analyse	<u>20- bis 40-fache Beschleunigung: Entscheidungsfindung innerhalb von Stunden statt Tagen</u>
<a href="#">ParStream</a>	Beschleunigung von Datenbanken und Datenanalyse mit Grafikprozessoren	<u>10-fache oder noch größere Beschleunigung: Milliarden von Datensätzen werden in Sekunden durchsucht</u>
<a href="#">Tanay Data Analytics von Fuzzy Logix</a>	Datenbank-integrierte Analytik-Engine	<u>Beschleunigung von Finanzsimulationen, Data-Mining und statistischen Verfahren mit Grafikprozessoren</u>
<a href="#">GPU-Quicksort</a>	Hoch optimierter Quicksort-Algorithmus	<u>10-fache Beschleunigung: Über 16 Millionen Gleitkommazahlen werden in weniger als 0,5 s sortiert</u>
<a href="#">AccelerEyes Jacket für MATLAB</a>	Nahtlose Beschleunigung von MATLAB mit Jacket™	<u>Bis zur 100-fachen Beschleunigung je nach Anwendung</u>
<a href="#">Wolfram Mathematica</a>	Analysen in symbolischer Mathematik (Mathematica)	Deutliche Leistungssteigerungen bei unterschiedlichen Anwendungen, unter anderem für lineare Algebra, Bildbearbeitung, Finanzsimulationen und Fourier-Transformationen

---

# 3. Der Algorithmus beim Sortieren

## Phase 1

- Eingeegebene Zahlen werden in gleich so unterteilt, sodass es gleich viele Blöcke wie Treads gibt
    - ermöglicht später synchronisiertes Arbeiten
    - Partitionierung
-

---

**Algorithm 2 GPU-Quicksort (First Phase GPU Kernel)**

---

```
function GQSORT(blocks, d,  $\hat{d}$ )
var global sstart, send, oldstart, oldend, blockcount
var block local lt, gt, pivot, start, end, s
var thread local i, lfrom, gfrom, lpivot, gpivot

(sstart→end, pivot, parent) ← blocksblockid ▷ Get the sequence block assigned to this thread block.
ltthreadid, gtthreadid ← 0, 0 ▷ Set thread local counters to zero.

i ← start + threadid ▷ Align thread accesses for coalesced reads.
for i < end, i ← i + threadcount do ▷ Go through the data...
  if si < pivot then ▷ counting elements that are smaller...
    ltthreadid ← ltthreadid + 1
  if si > pivot then ▷ or larger compared to the pivot.
    gtthreadid ← gtthreadid + 1

lt0, lt1, lt2, ..., ltsum ← 0, lt0, lt0 + lt1, ...,  $\sum_{i=0}^{threadcount} lt_i$  ▷ Calculate the cumulative sum.
gt0, gt1, gt2, ..., gtsum ← 0, gt0, gt0 + gt1, ...,  $\sum_{i=0}^{threadcount} gt_i$ 

if threadid = 0 then ▷ Allocate memory in the sequence this block is a part of.
  (seqsstart→send, oseqoldstart→oldend, blockcount) ← parent ▷ Get shared variables.
  lbeg ← FAA(sstart, ltsum) ▷ Atomic increment allocates memory to write to.
  gbeg ← FAA(send, -gtsum) - gtsum ▷ Atomic is necessary since multiple blocks access this
variable.

lfrom = lbeg + ltthreadid
gfrom = gbeg + gtthreadid

i ← start + threadid
for i < end, i ← i + threadcount do ▷ Go through data again writing elements
  if si < pivot then ▷ to the their correct position.
     $\neg s$ lfrom ← si ▷ If s is a sequence in d,  $\neg s$  denotes the corresponding
    lfrom ← lfrom + 1 ▷ sequence in  $\hat{d}$  (and vice versa).
  if si > pivot then
     $\neg s$ gfrom ← si
    gfrom ← gfrom + 1

if threadid = 0 then
  if FAA(blockcount, -1) = 0 then ▷ Check if this is the last block in the sequence to finish.
    for i ← sstart, i < send, i ← i + 1 do ▷ Fill in pivot value.
      di ← pivot
    lpivot ← median( $\neg seq$ oldstart,  $\neg seq$ (oldstart+sstart)/2,  $\neg seq$ sstart)
    gpivot ← median( $\neg seq$ send,  $\neg seq$ (send+oldend)/2,  $\neg seq$ oldend)
    result ← result ∪ {( $\neg seq$ oldstart→sstart, lpivot)}
    result ← result ∪ {( $\neg seq$ send→oldend, gpivot)}
```

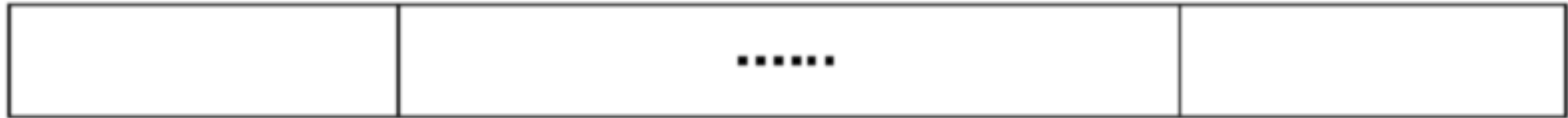
---

---

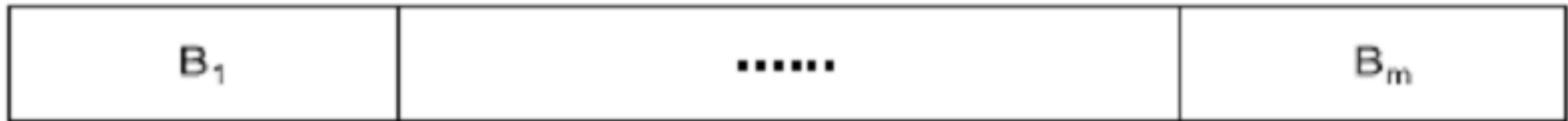
## Phase 2

- Jeder Thread bekommt nun seinen eigenen Block zugewiesen
    - Thread-Synchronisation nicht mehr nötig
  - Läuft ausschließlich auf dem GPU
  - Falls nötig, weiteres Partitionieren
-

# Partitionierung



(a) The sequence to be sorted is divided into  $m$  equally sized sections



(b) Thread blocks are assigned to the sections



(c) Each thread block goes through its assigned section...

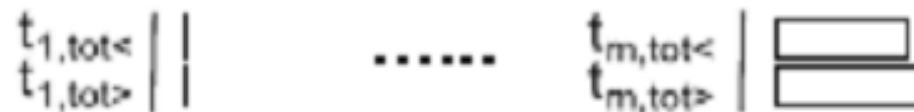
# Partitionierung (cont'd)



...keeping track of the number of elements above and below the pivot

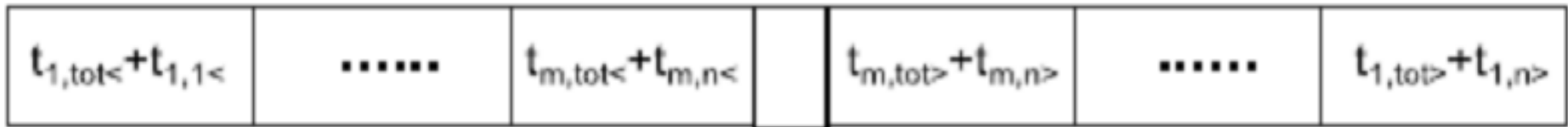


(d) Each thread block calculates the cumulative sum of the two types of elements seen



(e) A cumulative sum of each thread blocks total is calculated

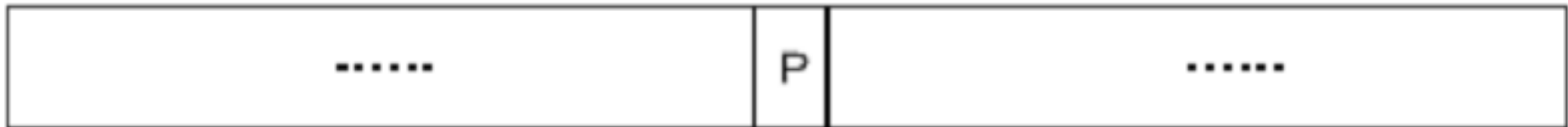
# Partitionierung (cont'd)



(f) Each thread uses the cumulative sums to find out where to write



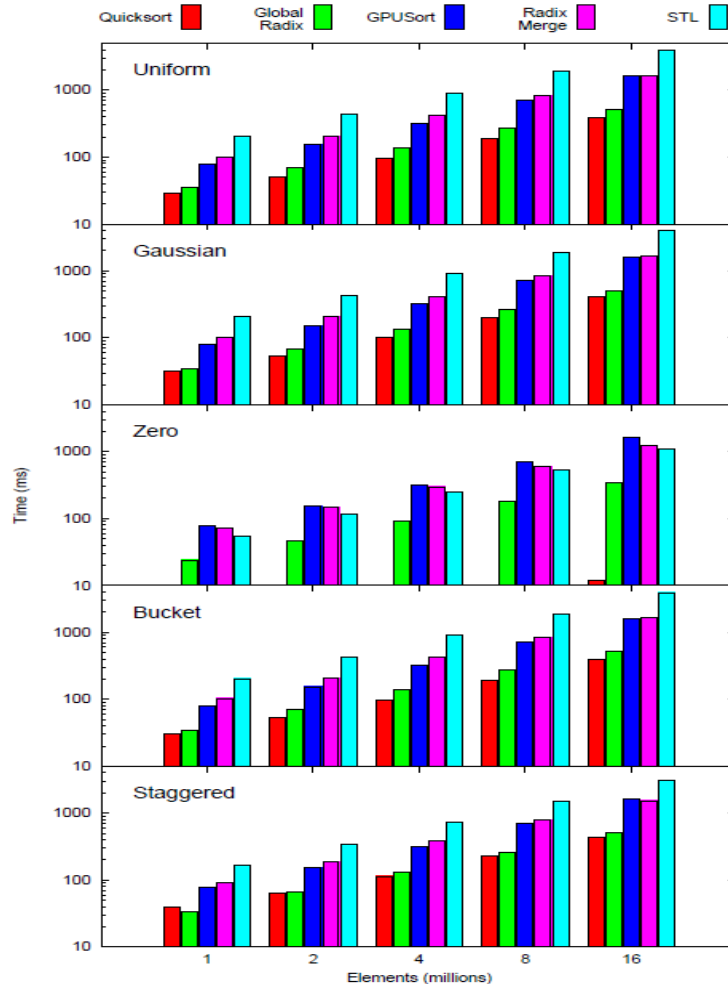
(g) Each thread block goes through its assigned data, writing it to the auxiliary array



(h) Block  $B_m$  fills the gap between the left and right subsequence with the pivot value

# 4. Beispiele

8800 GTX



= Standard Template Library

Zufällige Zahlen zwischen 0 und  $2^{31}$

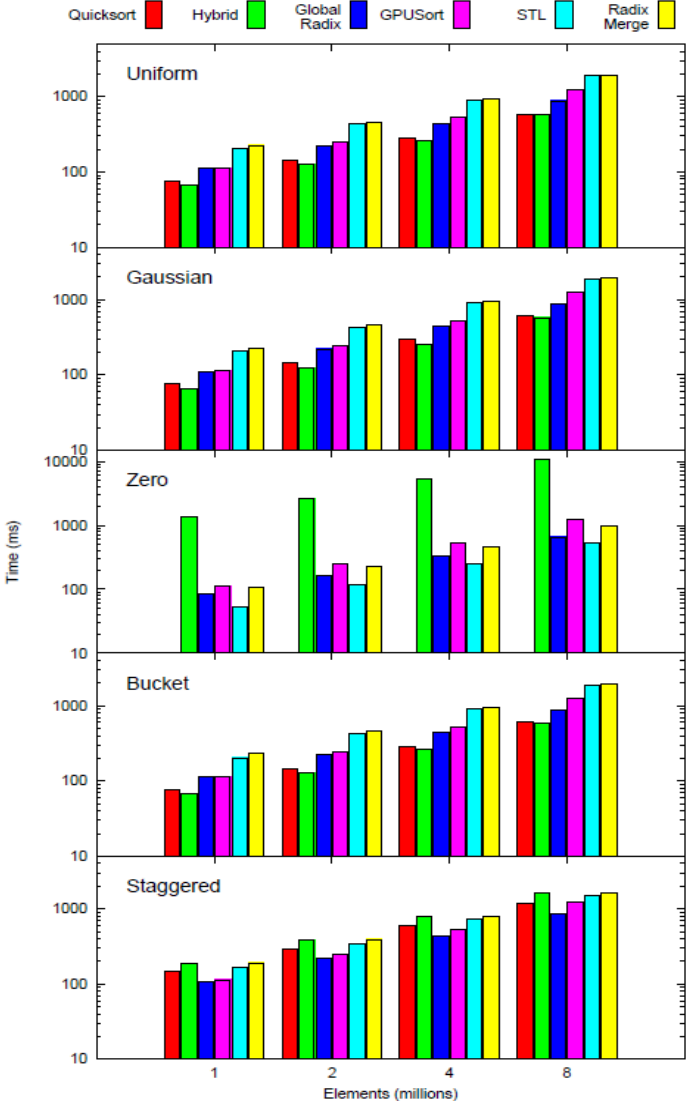
Durchschnitt von 4 zufällig gewählten Zahlen aus dem Uniform-Teil

Eine zufällige Konstante

Einteilung in Blöcke mit verschiedenen Werten



# 8600GTS



---

# Quellen

[http://www.nvidia.com/object/cuda\\_home\\_new.html](http://www.nvidia.com/object/cuda_home_new.html)

<http://www.cse.chalmers.se/research/group/dcs/gpuqsortdcs.html>

---