

Einführung in die Technische Informatik

GPU Programmierung

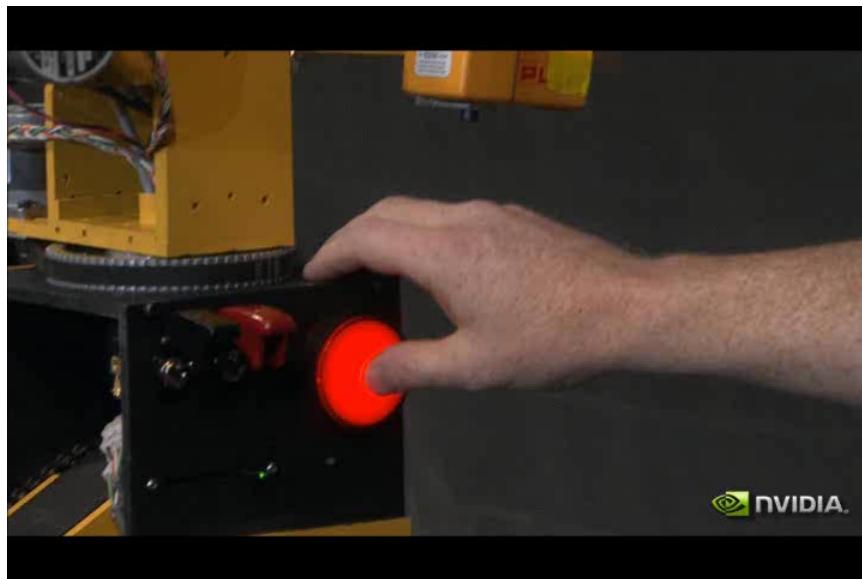
Wolfgang E. Nagel (wolfgang.nagel@tu-dresden.de)

Guido Juckeland (guido.juckeland@tu-dresden.de)



WARUM PARALLEL RECHNEN?

Eine Recheneinheit – Vielseitig, aber alleine



www.nvidia.com/content/nvision2008/art_science/index.html



TECHNISCHE
UNIVERSITÄT
DRESDEN

Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 3



Viele Recheneinheiten – Dumm, aber stark im Team



www.nvidia.com/content/nvision2008/art_science/index.html

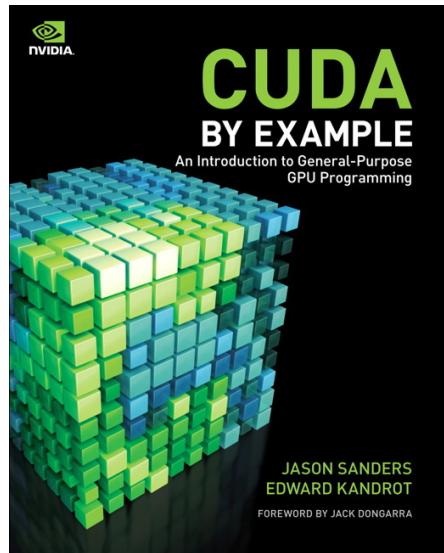
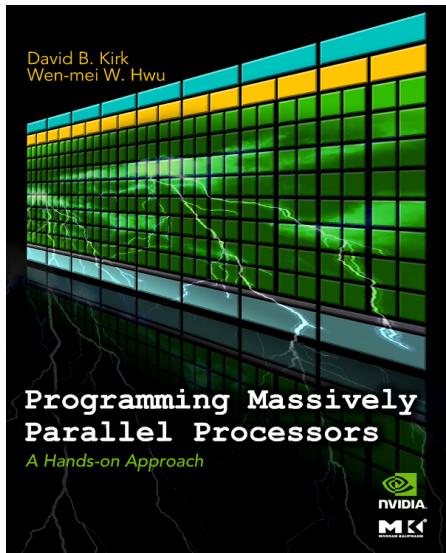


TECHNISCHE
UNIVERSITÄT
DRESDEN

Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 4



Literatur



Beide Bücher liegen in der SLUB in ausreichender Zahl zur Ausleihe bereit



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 5



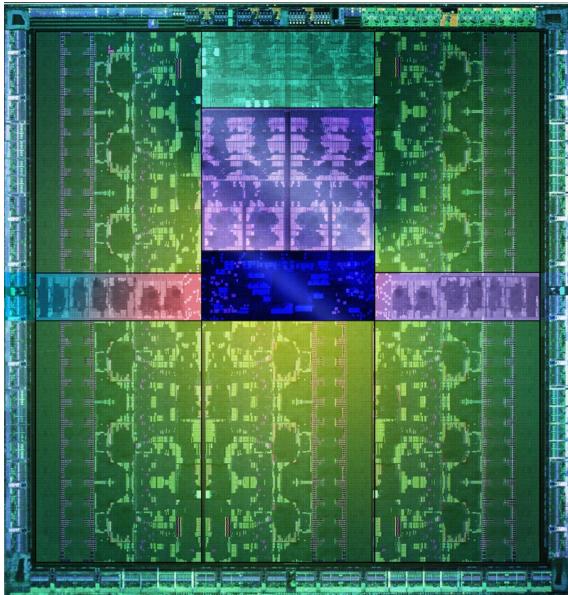
GRAFIKKARTEN ALS VIELKERNPROZESSOREN



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 6



NVIDIA Kepler GK110 Chip



- 2880 CUDA Kerne
- 15 SMX Mehrkerneinheiten
- 384-bit Memory Controller
- Bis zu 24GB of GDDR5 Speicher
- Zweite Generation ECC
- Mindestens 1.5 TFLOPS DP FP64
- Ziel: 250 GB/s Bandbreite

vr-zone.com/articles/nvidia-s-monster-gpu-for-tesla-k20-2013-geforce-and-quadro-cards/15884.html



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 7



NVIDIA Kepler GK110 Aufbau



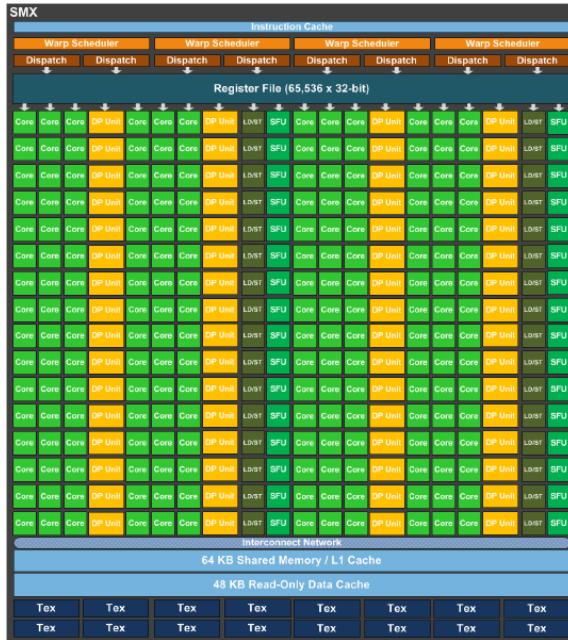
www.anandtech.com/show/5840/gtc-2012-part-1-nvidia-announces-gk104-based-tesla-k10-gk110-based-tesla-k20/2



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 8



NVIDIA Kepler GK110 Mehrkerneinheit



SMX
„Streaming
Multiprocessor“

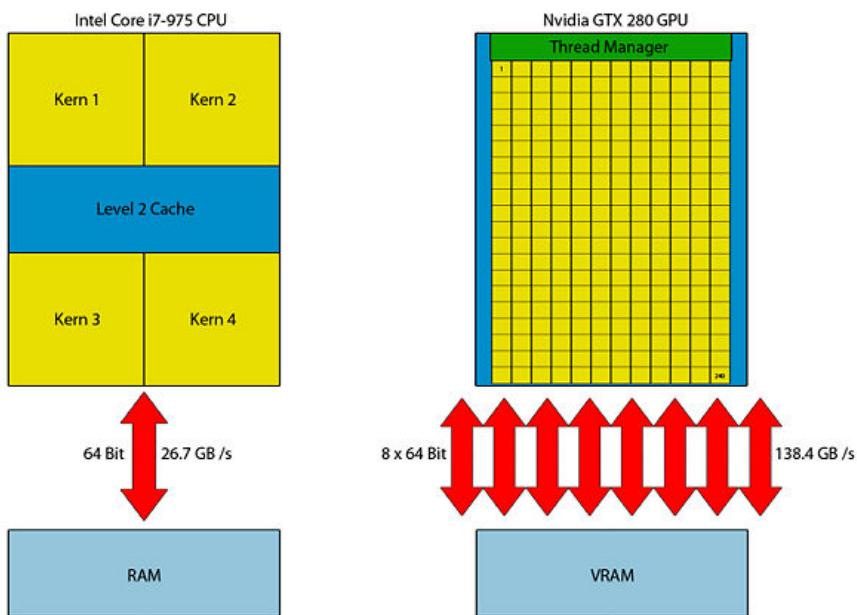
www.anandtech.com/show/5840/gtc-2012-part-1-nvidia-announces-gk104-based-tesla-k10-gk110-based-tesla-k20/2



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 9



Der kleine Unterschied



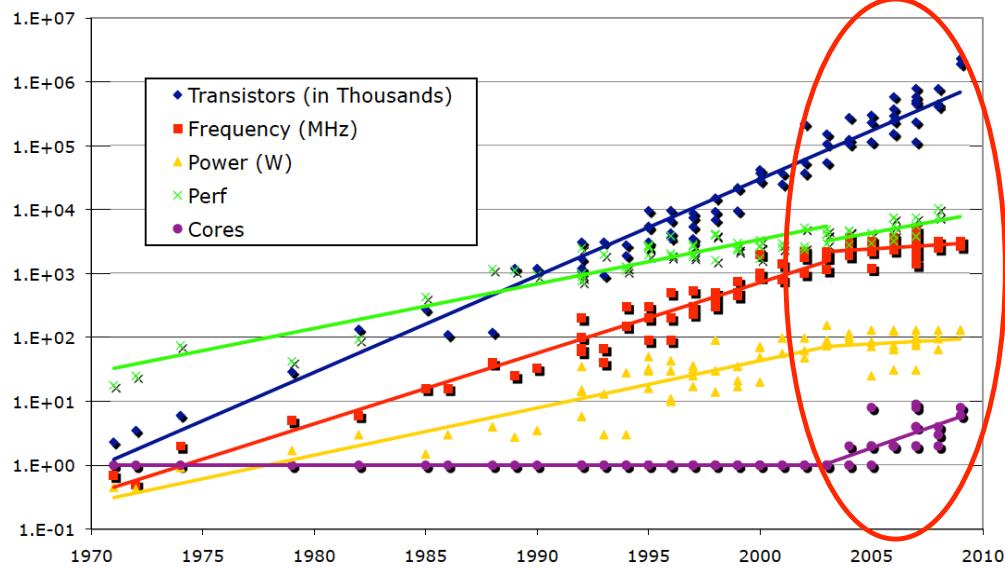
winfwiki.wi-fom.de/index.php/Bild:GPGPU_CPU_vs_GPU.jpg



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 10



Nur mit mehr Rechenkernen gibt's mehr Leistung



Kathy Yelick: „Ten Ways to Waste a Parallel Computer“

Keynote ISCA 2009. The 36th International Symposium on Computer Architecture
(mit Daten von Kunle Olukotun, Lance Hammond, Herb Sutter, Burton Smith, Chris Batten und Krste Asanović)



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 11



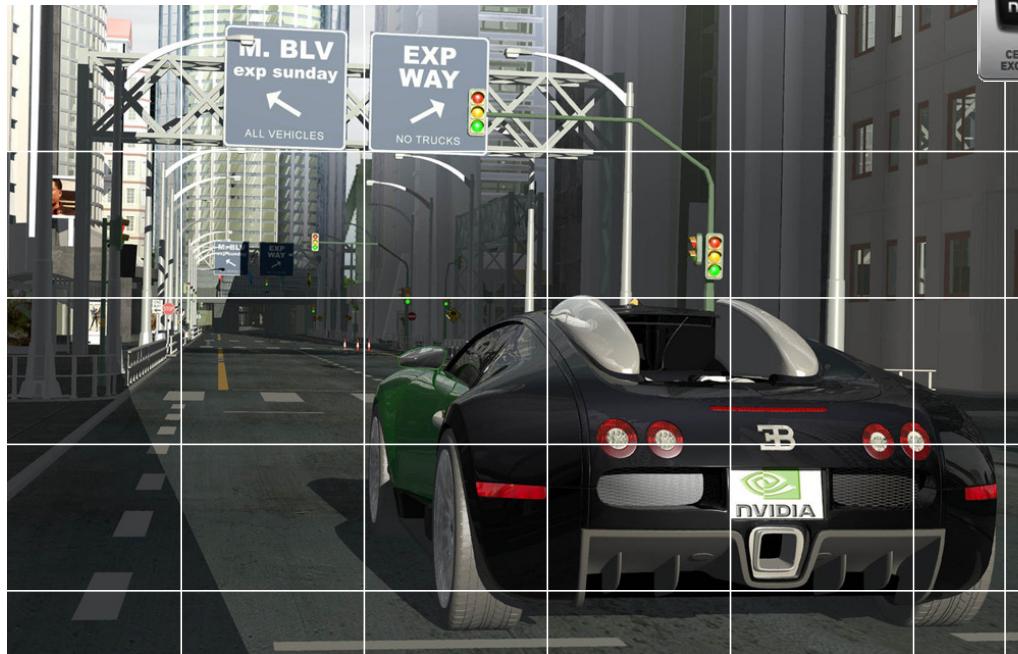
GRAFIKKARTEN ZUM RECHNEN NUTZEN



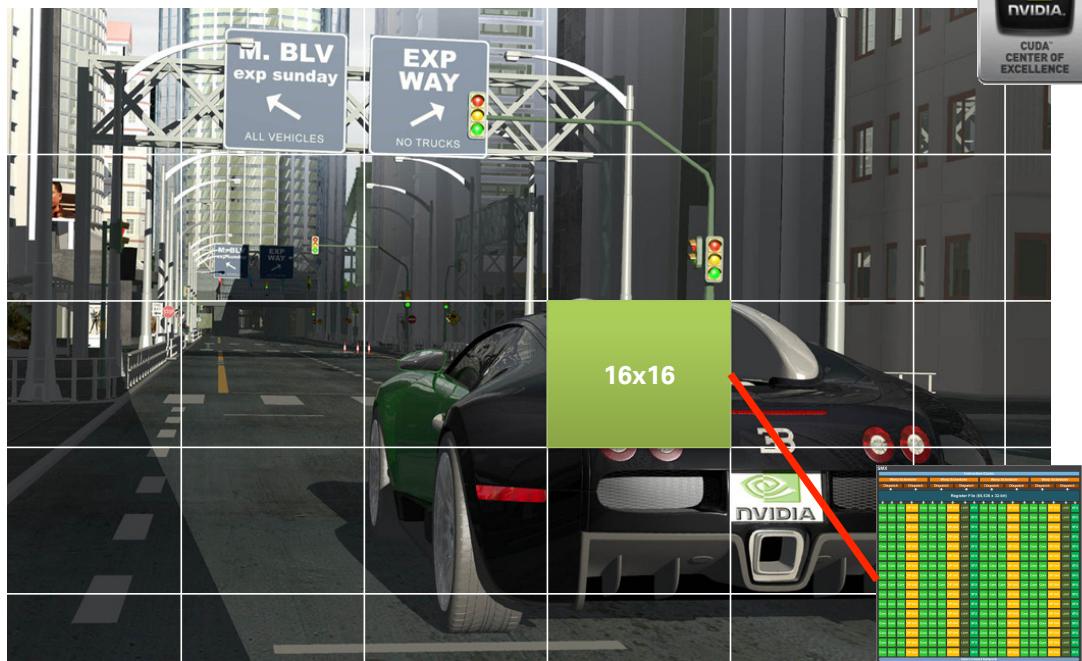
Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 12



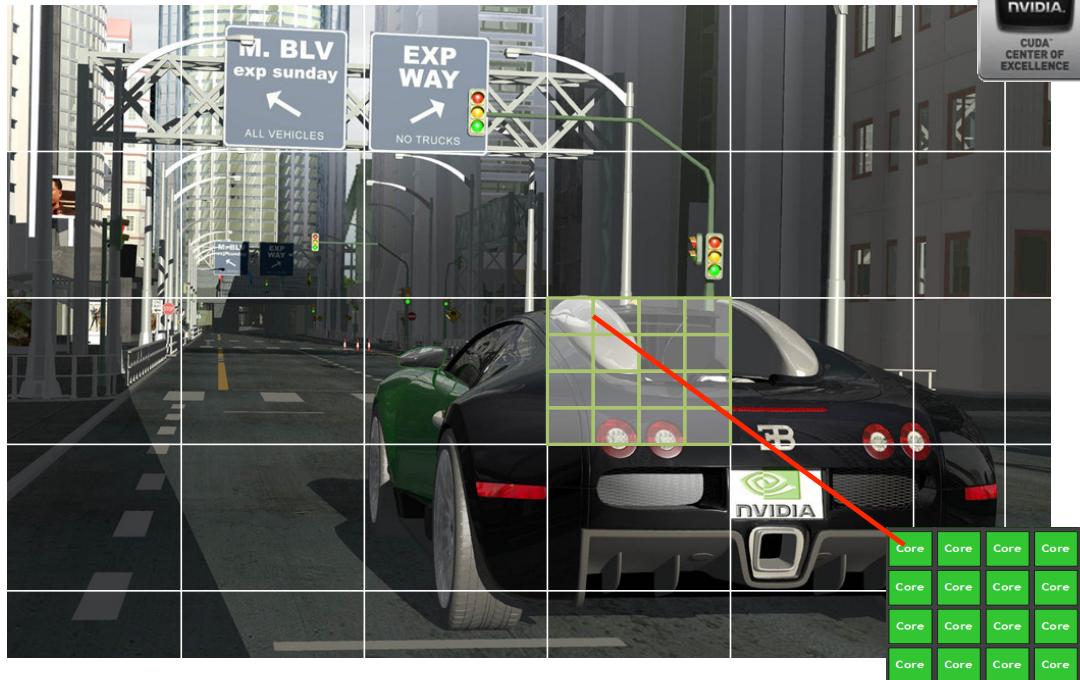
Zerlegung eines Bildes in Gitterblöcke



Block wird SMX zugewiesen



Untergitter (Pixel) werden Cores zugewiesen



Verschiedene Universen (Beispiele)



Spracherweiterungen

- CUDA
- OpenCL

Verschiedene Universen (Beispiele)



Bibliotheken

- cuBLAS
- cuFFT
- Thrust

Grafik

- OpenGL
- DirectX

Prozessmodellierung

- LabView
- Simulink



Mathematikprogramme

- MATLAB
- Mathematica
- IDL

Spracherweiterungen

- CUDA
- OpenCL
- DirectCompute

PRAGMA-Ansätze

- OpenACC
- OpenHMPP
- OpenMP

Interpretersprachen

- Python: PyCUDA
- Perl: KappaCUDA

Hochsprachen

- C/C++
- FORTRAN
- JAVA (RootBeer, JCUDA)
- Erlang

Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 17



CUDA Community Showcase



VexCL				
<input type="text" value="Suchen"/> <input type="button" value="Nach Veröffentlichung..."/> <input type="button" value="Projekt einreichen"/> <input type="checkbox"/> Nach Art der Anwendung filtern <input type="checkbox"/> Nach Art des Inhalts filtern <input type="checkbox"/> Nach Art der Organisation filtern <input type="checkbox"/> Numerische... <input type="checkbox"/> Bildgebung <input type="checkbox"/> Anwendung <input type="checkbox"/> Präsentation <input type="checkbox"/> DCC (Digital... <input type="checkbox"/> Medizinsch... <input type="checkbox"/> Wissenschaft <input type="checkbox"/> Multimedia <input type="checkbox"/> EDA (Electro... <input type="checkbox"/> Numerik <input type="checkbox"/> Signaverarb... <input type="checkbox"/> Code <input type="checkbox"/> Artikel <input type="checkbox"/> Finanzen <input type="checkbox"/> Life Sciences <input type="checkbox"/> Video & Audio <input type="checkbox"/> Bibliotheken <input type="checkbox"/> Sonstiges <input type="checkbox"/> Spielephysik <input type="checkbox"/> Öl & Gas <input type="checkbox"/> Hochschule <input type="checkbox"/> Unternehmen <input type="checkbox"/> Forschung				

www.nvidia.de/object/cuda_apps_flash_de.html#state=home



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 18



Supercomputer TOP 10

Rank	Site	System	Cores	Rmax (TFlop/s)	Rpeak (TFlop/s)	Power [kW]
1	DOE/SC/Oak Ridge National Laboratory United States	Titan - Cray XK7 , Opteron 6274 16C 2.200GHz, Cray Gemini interconnect, NVIDIA K20x Cray Inc.	560640	17590.0	27112.5	8209
2	DOE/NNSA/LLNL United States	Sequoia - BlueGene/Q, Power BQC 16C 1.60 GHz, Custom IBM	1572864	16324.8	20132.7	7890
3	RIKEN Advanced Institute for Computational Science (AICS) Japan	K computer, SPARC64 VIIIfx 2.0GHz,Tofu interconnect Fujitsu	705024	10510.0	11280.4	12660
4	DOE/SC/Argonne National Laboratory United States	Mira - BlueGene/Q, Power BQC 16C 1.60GHz, Custom IBM	786432	8162.4	10066.3	3945
5	Forschungszentrum Juelich (FZJ) Germany	JUQUEEN - BlueGene/Q, Power BQC 16C 1.600GHz, Custom Interconnect Custom	393216	4141.2	5033.2	1970
6	Leibniz Rechenzentrum Germany	SuperMUC - iDataPlex DX360M4, Xeon E5-2680 8C 2.70GHz, Infiniband FDR IBM	147456	2897.0	3185.1	3423
7	Texas Advanced Computing Center/Univ. of Texas United States	Stampede - PowerEdge C8220, Xeon E5-2680 8C 2.700GHz, Infiniband FDR, Intel Xeon Phi Dell	204900	2660.3	3959.0	
8	National Supercomputing Center in Tianjin China	Tianhe-1A - NUDTYH MPP, Xeon X5670 6C 2.93 GHz, NVIDIA 2050 NUDT	186368	2566.0	4701.0	4040
9	CINECA Italy	Fermi - BlueGene/Q, Power 5000, 713672912/11om IBM	163840	1725.5	2097.2	822
10	 Technische Universität Dresden Germany	DARPA Trial Subset - Power 775, POWER7 8C 3.836GHz, Custom Interconnect IBM	63360	1515.0	1944.4	3576

Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 19



Center for Information Services & High Performance Computing

Top500 11/13 - TOP 10

Rank	Installation Site	Manufacturer	Computer	Country	Cores	Rmax [Tflops]	Power [MW]
1	National Super Computer Center in Guangzhou	NUDT	Tianhe-2 (MilkyWay-2) TH-IVB-FEP Cluster, Intel Xeon E5-2692 12C 2.200GHz, Intel Xeon Phi	China	3120000	33863	17.8
2	DOE/SC/Oak Ridge National Laboratory	Cray Inc.	Titan Cray XK7 , Opteron 6274 16C 2.200GHz, NVIDIA K20x	USA	560640	17590	8.2
3	DOE/NNSA/LLNL	IBM	Sequoia BlueGene/Q, Power BQC 16C 1.60 GHz	USA	1572864	17173	7.9
4	RIKEN Advanced Institute for Computational Science (AICS)	Fujitsu	K computer SPARC64 VIIIfx 2.0GHz	Japan	705024	10510	12.7
5	DOE/SC/Argonne National Laboratory	IBM	Mira BlueGene/Q, Power BQC 16C 1.60GHz	USA	786432	8587	3.9
6	Swiss National Supercomputing Centre (CSCS)	Cray Inc.	Piz Daint Cray XC30, Xeon E5-2670 8C 2.600GHz, NVIDIA K20x	Switzerland	115984	6271	2.3
7	Texas Advanced Computing Center/Univ. of Texas	Dell	Stampede PowerEdge C8220, Xeon E5-2680 8C 2.700GHz, Intel Xeon Phi	USA	462462	5168	4.5
8	Forschungszentrum Juelich (FZJ)	IBM	JUQUEEN BlueGene/Q, Power BQC 16C 1.600GHz	Germany	393216	5009	2.3
9	DOE/NNSA/LLNL	IBM	Vulcan BlueGene/Q, Power BQC 16C 1.600GHz	USA	393216	4293	2.0
10	Leibniz Rechenzentrum	IBM	SuperMUC iDataPlex DX360M4, Xeon E5-2680 8C 2.70GHz	Germany	147456	2897	3.4



Center for Information Services & High Performance Computing

Green500 11/13 - TOP 10



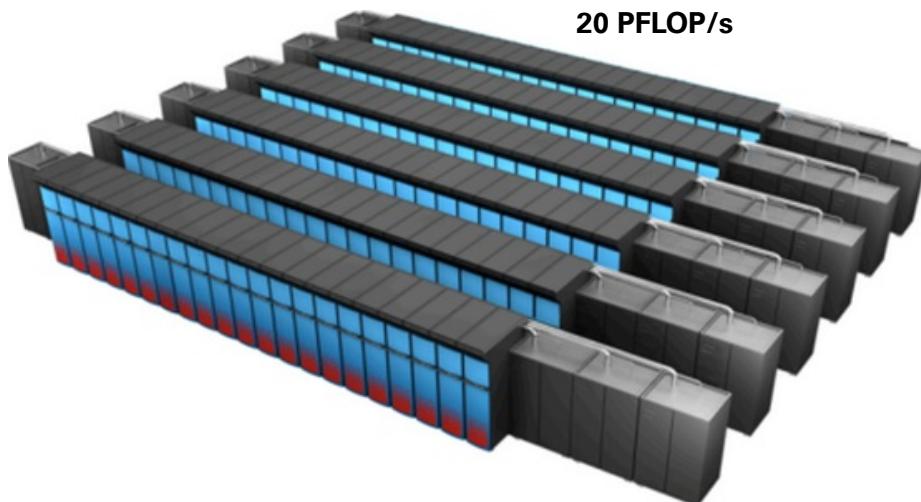
Rank	Installation Site	Manufacturer	Computer	Country	Cores	Rmax [Gflops]	MFLOPS/W	Top500 rank
1	GSIC Center, Tokyo Institute of Technology	NEC	TSUBAME-KFC LX 1U-4GPU/104Re-1G Cluster, Intel Xeon E5-2620v2 6C, 2.100GHz, NVIDIA K20x	Japan	2720	125	4503	311
2	Cambridge University	DELL	Wilkes Dell T620 Cluster, Intel Xeon E5-2630v2 6C 2.600GHz, NVIDIA K20	UK	5120	191	3632	166
3	Center for Computational Sciences, University of Tsukuba	Cray	HA-PACS TCA Cray 3623G4-SM Cluster, Intel Xeon E5-2680v2 10C 2.800GHz, NVIDIA K20x	Japan	4864	277	3518	134
4	Swiss National Supercomputing Centre (CSCS)	Cray	Piz Daint Cray XC30, Xeon E5-2670 8C 2.600GHz, NVIDIA K20x	Switzerland	115984	5587	3186	6
5	ROMEO HPC Center - Champagne-Ardenne	Bull SA	romeo Bull R421-E3 Cluster, Intel Xeon E5-2650v2 8C 2.600GHz, NVIDIA K20x	France	5720	255	3131	151
6	GSIC Center, Tokyo Institute of Technology	NEC/HP	TSUBAME 2.5 Cluster Platform SL390s G7, Xeon X5670 6C 2.930GHz, NVIDIA K20x	Japan	74358	2831	3069	11
7	University of Arizona	IBM	iDataPlex DX360M4, Intel Xeon E5-2650v2 8C 2.600GHz, NVIDIA K20x	USA	3080	145	2702	336
8	Max-Planck-Gesellschaft MPI/IPP	IBM	iDataPlex DX360M4, Intel Xeon E5-2680v2 10C 2.800GHz, NVIDIA K20x	Germany	15840	710	2629	49
9	Financial Institution	IBM	iDataPlex DX360M4, Intel Xeon E5-2680v2 10C 2.800GHz, NVIDIA K20x	USA	3264	146	2629	328
10	CSIRO	Xenon Systems	CSIRO GPU Cluster Nitro G16 3GPU, Xeon E5-2650 8C 2.000GHz, Nvidia K20m	Australia	4620	168	2359	260

TITAN at Oak Ridge National Lab

18.000 KEPLER GPUs



20 PFLOP/s



perilsofparallel.blogspot.de/2012/01/20-pflops-vs-10s-of-mloc-oak-ridge.html



ERSTE SCHRITTE CUDA INSTALLIEREN



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 23



developer.nvidia.com/category/zone/cuda-zone

The screenshot shows the NVIDIA Developer Zone homepage with a sidebar for the CUDA Zone. The CUDA Zone sidebar includes links for CUDA Overview, CUDA Toolkit, CUDA Samples, CUDA Documentation, Tool & Examples, and Benchmark/Tuning. Below this, there are sections for CUDA PRE-PRODUCTION, GET STARTED - PARALLEL COMPUTING, CUDA TOOLKIT, CUDA DOWNLOADS (highlighted with a red box), CUDA TOOLS & ECOSYSTEM, OPENACC, CUDA EDUCATION & TRAINING, and NVIDIA GPU COMPUTING DOCUMENTATION.

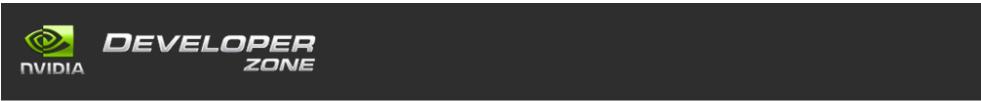
The screenshot shows the CUDA DOWNLOADS section for CUDA Toolkit 4.2. It features three main download paths: 1. Download Toolkit (64bit, 32bit), 2. Download Drivers (Win7/Vista Desktop, Win7/Vista Notebook, Win XP Desktop, 64bit, 32bit), and 3. Download SDK (64bit, 32bit). Below these are sections for CUDA 4.2 FOR LINUX (Fedora 14, Redhat 5.5, Ubuntu 10.04, OpenSUSE 11.2, SUSE Server 11 SP1) and CUDA 4.2 FOR MAC (Download Toolkit, Download Drivers (updated 06.11.12), Download SDK).

Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 24

ZIH
Center for Information Services &
High Performance Computing

Getting Started (Guides + Beispiele)

developer.nvidia.com/cuda/nvidia-gpu-computing-documentation



Home

NVIDIA GPU COMPUTING DOCUMENTATION

Find all the latest documentation for the CUDA Toolkit, Libraries, Tools and SDK's below.

Individual SDK code samples can be found [here](#).

Documentation of legacy versions of the CUDA Toolkit and SDKs are included in the installable legacy downloads available on the [CUDA Toolkit Archive Page](#).

We recommend that all developers upgrade to the latest version, download the [Latest Version of the CUDA Toolkit](#) today.

CUDA Getting Started Guide (Windows)

[Download](#)

This guide will show you how to install and check the correct operation of the CUDA development tools in Windows.

CUDA Getting Started Guide (Linux)

[Download](#)

This guide will show you how to install and check the correct operation of the CUDA development tools in Linux.

CUDA Getting Started Guide (Mac OS X)

[Download](#)

This guide will show you how to install and check the correct operation of the CUDA development tools in Mac OS X.

Getting Started with CUDA SDK samples

[Download](#)

This guide covers the introductory CUDA SDK samples beginning CUDA developers should review before developing your own projects.



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 25



Schrittweise zur erfolgreichen CUDA-Installation

Installationsschritte

- Habe ich Zugang zu CUDA-fähiger Hardware?



GeForce Desktop Products	
GPU	Compute Capability
GeForce GTX 690	3.0
GeForce GTX 680	3.0
GeForce GTX 670	3.0
GeForce GTX 660 Ti	3.0
GeForce GTX 660	3.0
GeForce GTX 650	3.0
GeForce GTX 560 Ti	2.1
GeForce GTX 550 Ti	2.1
GeForce GTX 460	2.1
GeForce GTS 450	2.1
GeForce GTS 450*	2.1
GeForce GTX 590	2.0
GeForce GTX 580	2.0
GeForce GTX 570	2.0



CUDA GPUS	
NVIDIA GPUs power millions of desktops, notebooks, workstations and supercomputers around the world, accelerating computationally-intensive tasks for consumers, professionals, scientists, and researchers.	
Find out all about CUDA and GPU Computing by attending our GPU Computing Webinars and joining our free-to-join CUDA Registered developer Program .	
• Learn about Tesla for technical and scientific computing	
• Learn about Quadro for professional visualization	

CUDA-Enabled Tesla GPU Computing Products

Tesla Workstation Products		Tesla Data Center Products	
GPU	Compute Capability	GPU	Compute Capability
Tesla C2075	2.0	Tesla K20	3.5
Tesla C2080/C2070	2.0	Tesla K10	3.0
Tesla C1060	1.3	Tesla M2050/M2070/M2075/M2090	2.0
Tesla C870	1.0	Tesla S1070	1.3
Tesla D870	1.0	Tesla M1060	1.3
		Tesla S870	1.0

CUDA-Enabled Quadro Products

Quadro Desktop Products		Quadro Mobile Products	
GPU	Compute Capability	GPU	Compute Capability
Quadro K5000	3.0	Quadro K500M	3.0
Quadro 4000	2.0	Quadro 5010M	2.0
Quadro 5000	2.0	Quadro 5000M	2.0
Quadro 4000	2.0	Quadro 4000M	2.1
Quadro 4000 for Mac	2.0	Quadro 3000M	2.1



Schrittweise zur erfolgreichen CUDA-Installation



Installationsschritte

- Habe ich Zugang zu CUDA-fähiger Hardware?
- Grafikkartentreiber mit CUDA-Fähigkeiten herunterladen



developer.nvidia.com/cuda/cuda-downloads

The screenshot shows the CUDA DOWNLOADS section of the NVIDIA Developer Zone. It features a prominent 'CUDA TOOLKIT 4.2 CURRENT PRODUCTION RELEASE' banner. Below it, there are three main sections: 'CUDA 4.2 FOR WINDOWS', 'CUDA 4.2 FOR LINUX', and 'CUDA 4.2 FOR MAC'. Each section contains numbered steps (1, 2, 3) leading to specific download links. Step 1 is 'Download Toolkit', Step 2 is 'Download Drivers (v 301.32 or 301.27)', and Step 3 is 'Download SDK'. The Windows section includes links for Win7/Vista Desktop (64bit, 32bit), Win7/Vista Notebook (64bit, 32bit), and Win XP Desktop (64bit, 32bit). The Linux section includes links for Fedora 14 (64bit, 32bit), Redhat 5.5 (64bit, 32bit), Ubuntu 10.04 (64bit, 32bit), OpenSUSE 11.2 (64bit, 32bit), and SUSE Server 11 SP1 (64bit, 32bit). The Mac section includes links for Download Toolkit (DOWNLOAD), Download Drivers (updated 06.11.12) (DOWNLOAD), and Download SDK (DOWNLOAD).



Schrittweise zur erfolgreichen CUDA-Installation



Installationsschritte

- Habe ich Zugang zu CUDA-fähiger Hardware?
- Grafikkartentreiber mit CUDA-Fähigkeiten herunterladen
- Prüfen, ob das System CUDA-Grafikkartentreiber unterstützt
(vor allem bei LINUX wichtig!)



Schrittweise zur erfolgreichen CUDA-Installation



Installationsschritte

- Habe ich Zugang zu CUDA-fähiger Hardware?
- Grafikkartentreiber mit CUDA-Fähigkeiten herunterladen
- Prüfen, ob das System CUDA-Grafikkartentreiber unterstützt
(vor allem bei LINUX wichtig!)
- Prüfen, ob ein Compiler installiert ist, der CUDA unterstützt
(LINUX: gcc, Windows: MS Visual Studio mit C++, Mac OS X: gcc)

Schrittweise zur erfolgreichen CUDA-Installation



Installationsschritte

- Habe ich Zugang zu CUDA-fähiger Hardware?
- Grafikkartentreiber mit CUDA-Fähigkeiten herunterladen
- Prüfen, ob das System CUDA-Grafikkartentreiber unterstützt
(vor allem bei LINUX wichtig!)
- Prüfen, ob ein Compiler installiert ist, der CUDA unterstützt
(LINUX: gcc, Windows: MS Visual Studio mit C++, Mac OS X: gcc)
- CUDA-fähigen Grafikkartentreiber installieren

Schrittweise zur erfolgreichen CUDA-Installation



Installationsschritte

- Habe ich Zugang zu CUDA-fähiger Hardware?
- Grafikkartentreiber mit CUDA-Fähigkeiten herunterladen
- Prüfen, ob das System CUDA-Grafikkartentreiber unterstützt
(vor allem bei LINUX wichtig!)
- Prüfen, ob ein Compiler installiert ist, der CUDA unterstützt
(LINUX: gcc, Windows: MS Visual Studio mit C++, Mac OS X: gcc)
- CUDA-fähigen Grafikkartentreiber installieren
- CUDA Toolkit herunterladen und installieren

Schrittweise zur erfolgreichen CUDA-Installation



Installationsschritte

- Habe ich Zugang zu CUDA-fähiger Hardware?
- Grafikkartentreiber mit CUDA-Fähigkeiten herunterladen
- Prüfen, ob das System CUDA-Grafikkartentreiber unterstützt
(vor allem bei LINUX wichtig!)
- Prüfen, ob ein Compiler installiert ist, der CUDA unterstützt
(LINUX: gcc, Windows: MS Visual Studio mit C++, Mac OS X: gcc)
- CUDA-fähigen Grafikkartentreiber installieren
- CUDA Toolkit herunterladen und installieren
- Zusätzlich kann man den GPU Computing SDK herunterladen und installieren und sich Beispielprogramme anschauen



ERSTE SCHRITTE CUDA NUTZEN



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 33



Mein erstes CUDA-Programm!

vec2zero.cu

```
#define NELEMENTS 16;           //16 Elemente
int main( int argc, char *argv[] ) {

    float veca[NELEMENTS];      //Deklaration Vektor veca
    int i;                      //Laufvariable für Schleife

    for ( i=0; i<NELEMENTS; i++ ) { //Schleife über Elemente
        veca[i] = 0.0f;          //Zuweisung der Elemente
    }

    return 0;
}
```



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 34



Compilieren (LINUX)



```
bussmann@gpunode> nvcc vec2zero.cu -o vec2zero  
bussmann@gpunode> ./vec2zero
```

Das funktioniert nur bei korrekter Softwareinstallation und wenn der Code auf einem CUDA-fähigen Device ausgeführt werden kann.

Ohne GPU erhält man den Fehler:

```
NVIDIA: could not open the device file /dev/nvidiactl (No such device or address).
```



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 35



Die GPU wirklich nutzen: Einen „Kernel“ schreiben

```
#define NELEMENTS 16;  
__global__ void kernel( void ) { //ein Kernel, der nix macht  
}  
  
int main( int argc, char *argv[] ) {  
  
    float veca[NELEMENTS];  
    int i;  
  
    for ( i=0; i<NELEMENTS; i++ ) {  
        veca[i] = 0.0f;  
    }  
  
    return 0;  
}
```



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 36



Die GPU wirklich nutzen: Einen „Kernel“ aufrufen



```
#define NELEMENTS 16;  
  
__global__ void kernel( void ) { //ein Kernel, der nix macht  
}  
  
int main( int argc, char *argv[] ) {  
  
    float veca[NELEMENTS];  
    int i;  
  
    kernel<<<1,1>>>(); //Tut nix, will nur spielen  
    for ( i=0; i<NELEMENTS; i++ ) {  
        veca[i] = 0.0f;  
    }  
  
    return 0;  
}
```



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 37



Zwischenspiel: HOST und DEVICE



„HOST“



„DEVICE“



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 38



Speicher auf GPU reservieren und freigeben



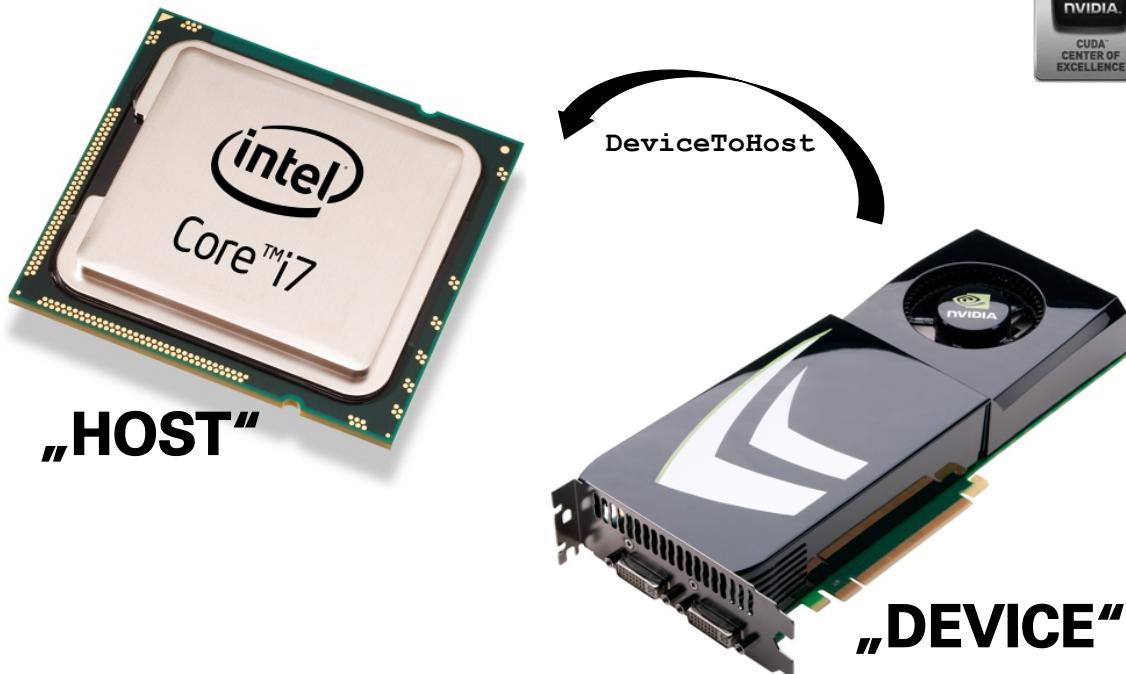
```
int main( int argc, char *argv[] ) {  
  
    float *devicepointer;                                //Zeiger auf float  
  
    //Die folgende Zeile reserviert Speicher für 1 float auf der GPU  
    //und setzt den Pointer veca_elem_dev auf diesen Speicherbereich  
    cudaMalloc( (void**)&devicepointer, sizeof(float) );  
    //Die folgende Zeile gibt den GPU-Speicher, auf den veca_elem_dev  
    //zeigt auf der GPU frei.  
    cudaFree( devicepointer );  
  
    return 0;  
}
```



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 39



Datentransfer: cudaMemcpy



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 40



Speicher von DEVICE nach HOST kopieren



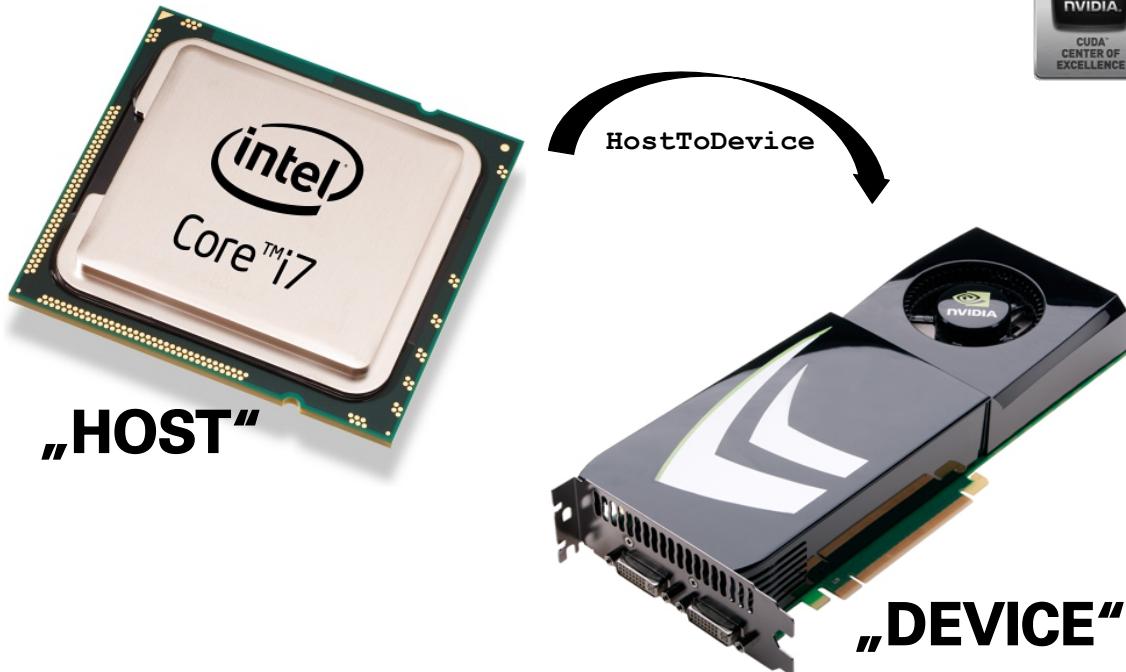
```
#define NELEMENTS 16;

int main( int argc, char *argv[] ) {

    float hostvariable;          //float-Variablen auf HOST
    float *devicepointer;        //Zeiger auf Speicher auf DEVICE

    cudaMemcpy( &hostvariable,      //Zeiger auf HOST-Variablen
               devicepointer,       //Zeiger auf DEVICE-Speicher
               sizeof(float),       //Größe des Speichers: 1 float
               cudaMemcpyDeviceToHost );
    return 0;
}
```

Datentransfer: cudaMemcpy



Speicher von HOST nach DEVICE kopieren



```
#define NELEMENTS 16;

int main( int argc, char *argv[] ) {

    float hostvariable;           //float-Variable auf HOST
    float *devicepointer;         //Zeiger auf Speicher auf DEVICE

    cudaMemcpy( devicepointer,      //Zeiger auf DEVICE-Speicher
                &hostvariable,      //Zeiger auf HOST-Variable
                sizeof(float),       //Größe des Speichers: 1 float
                cudaMemcpyHostToDevice );
}

return 0;
}
```



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 43



Ein kleiner Blick in die Familie cudaMemcpy



cudaMemcpy	cudaMemcpyToArray
cudaMemcpy2D	cudaMemcpyAsync
cudaMemcpy2DArrayToArray	cudaMemcpyFromArray
cudaMemcpy2DAsync	cudaMemcpyFromArrayAsync
cudaMemcpy2DFromArray	cudaMemcpyFromSymbol
cudaMemcpy2DFromArrayAsync	cudaMemcpyFromSymbolAsync
cudaMemcpy2DToArray	cudaMemcpyPeer
cudaMemcpy2DToArrayAsync	cudaMemcpyPeerAsync
cudaMemcpy3D	cudaMemcpyToArray
cudaMemcpy3DAsync	cudaMemcpyToArrayAsync
cudaMemcpy3DPeer	cudaMemcpyToSymbol
cudaMemcpy3DPeerAsync	cudaMemcpyToSymbolAsync



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 44



Fehler abfangen mit `cudaError_t`



CUDA-Funktionen geben einen Wert vom Typ `cudaError_t` zurück:

```
typedef enum cudaError cudaError_t
```

Dieser kann verschiedene Werte annehmen:

<code>cudaSuccess</code>	The API call returned with no errors. In the case of query calls, this can also mean that the operation being queried is complete (see <code>cudaEventQuery()</code> and <code>cudaStreamQuery()</code>).
<code>cudaErrorMissingConfiguration</code>	The device function being invoked (usually via <code>cudaLaunch()</code>) was not previously configured via the <code>cudaConfigureCall()</code> function.
<code>cudaErrorMemoryAllocation</code>	The API call failed because it was unable to allocate enough memory to perform the requested operation.
<code>cudaErrorInitializationError</code>	The API call failed because the CUDA driver and runtime could not be initialized.
<code>cudaErrorLaunchFailure</code>	An exception occurred on the device while executing a kernel. Common causes include dereferencing an invalid device pointer and accessing out of bounds shared memory. The device cannot be used until <code>cudaThreadExit()</code> is called. All existing device memory allocations are invalid and must be reconstructed if the program is to continue using CUDA.
...	

Den Fehler interpretieren

Zur Interpretation von `cudaError_t` kann man die Funktion



```
const char* cudaGetString( cudaError_t error)
```

nutzen:

```
cudaError_t error = cudaGetLastError();
printf( "CUDA error: %s\n", cudaGetString( error ) );
```

Fehler abfangen mit cutilXXXX

Der CUDA-SDK bietet einige Funktionen zum Fehlerabfangen



```
cutilSafeCall( ... );           //prüft Funktionsrückgabe  
                                //für CUDA-Funktionen  
  
cutilSafeCall( cudaMalloc(...) );  
  
cutilCheckMsg( ... );          //prüft kernel-Aufrufe  
kernel<<<1,1>>>();  
cutilCheckMsg("kernel failed\n");  
  
cutilCheckError( ... );        //prüft SDK-Aufrufe
```



TECHNISCHE
UNIVERSITÄT
DRESDEN

Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 47



Fehler abfangen mit HandleError (CUDA by Example)



```
static void HandleError( cudaError_t err, const char *file, int line ) {  
  
    if (err != cudaSuccess) {  
        printf( "%s in %s at line %d\n", cudaGetErrorString( err ), file, line );  
        exit( EXIT_FAILURE );  
    }  
}  
  
#define HANDLE_ERROR( err ) (HandleError( err, __FILE__, __LINE__ ))
```



TECHNISCHE
UNIVERSITÄT
DRESDEN

Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 48



Ein richtiges CUDA-Programm



```
#define NELEMENTS 16;  
__global__ void settozero( float *elem ) {  
    *elem = 0.0f;  
}  
  
int main( int argc, char *argv[] ) {  
  
    float veca[NELEMENTS];  
    float *devicemem;  
    int i;  
  
    cudaMalloc( (void**)&devicemem, sizeof(float) );  
    for ( i=0; i<NELEMENTS; i++ ) {  
        settozero<<<1,1>>>( devicemem );  
        cudaMemcpy( &veca[i], devicemem, sizeof(float), cudaMemcpyDeviceToHost );  
    }  
  
    cudaFree( devicemem );  
    return 0;  
}
```



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 49



ZUSAMMENFASSUNG



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 50



Kernel und Speicherverwaltung



Kernel werden auf der GPU (dem Device) aufgerufen. Sie können nur auf Speicher auf der GPU zugreifen. Sie verhalten sich wie C-Funktionen.

```
__global__ returntype kernelname( paramdec1, paramdec2, ... ) { ... }

kernelname<<<dim3 grid, dim3 block>>>( param1, param2, ... );
```

Speicherverwaltung:

```
cudaError_t cudaMalloc( void **devPtr, size_t size )
cudaError_t cudaFree( void *devPtr)
cudaError_t cudaMemcpy ( void *destination,           //Ziel
                      const void *source,          //Quelle
                      size_t count,               //Größe
                      enum cudaMemcpyKind direction ) //Richtung

enum cudaMemcpyKind: cudaMemcpyHostToDevice,
                     cudaMemcpyDeviceToHost, ...
```

Fehlerbehandlung



Fehler haben den Typ

`cudaError_t`

Sie werden entweder direkt von einer Funktion zurückgegeben oder können mit Hilfe der Funktion

```
cudaError_t cudaGetLastError( void )
bestimmt werden.
```

Der zugehörige Wert kann mit Hilfe der Funktion

```
const char* cudaGetString( cudaError_t error )
in einen Fehler-String konvertiert werden.
```



KERNEL UND THREADS

Kernel-Aufrufe



Ein **Kernel** ist ein Stück Code, das für die Ausführung auf der GPU übersetzt wird.
Kernel können auf allen CUDA-fähigen GPUs einzeln und hintereinander vom Host aufgerufen werden:

```
int main( int argc, char *argv[] ) {  
  
    kernel1<<<..., ...>>>(...);  
    kernel2<<<..., ...>>>(...);  
    kernel3<<<..., ...>>>(...);  
  
    return 0;  
}
```

Dabei wird gewartet, bis jeder Kernelaufruf abgearbeitet wurde, bevor der nächste Kernel aufgerufen wird. Die Abfolge der Kernelaufrufe bleibt somit strikt seriell.

Wie kann ich dann aber parallel arbeiten?



CUDA ermöglicht die parallele Ausführung von Kernel-Code. Dabei wird der Code eines Kernels parallel auf vielen Prozessoren gleichzeitig ausgeführt.

Jede Instanz eines solchen ausgeführten Kernels wird **Thread** genannt

```
int main( int argc, char *argv[] ) {  
  
    kernel1<<<...,>>>(...); //erzeugt viele, parallele Threads  
    kernel2<<<...,>>>(...); //erzeugt viele, parallele Threads  
    kernel3<<<...,>>>(...); //erzeugt viele, parallele Threads  
    return 0;  
}
```

Threads können, in Abhängigkeit der Daten, unterschiedliche Anweisungen ausführen. Um jedoch die maximale Leistung aus der GPU zu holen, sollte man immer im **SIMD-Modell** denken.

Exkurs: SIMD



SIMD steht für **S**ingle **I**nstruction, **M**ultiple **D**ata. Man kann dies so zusammenfassen, dass unabhängig von den Daten jeder Thread, d.h. jeder parallele Kernel-Aufruf, dieselben Instruktionen mit jeweils anderen Daten ausführt:

```
__global__ void settozero( float *elem ) {  
    *elem = 0.0f;  
}  
  
...  
  
settozero<<<...,>>> ( devicemem );
```

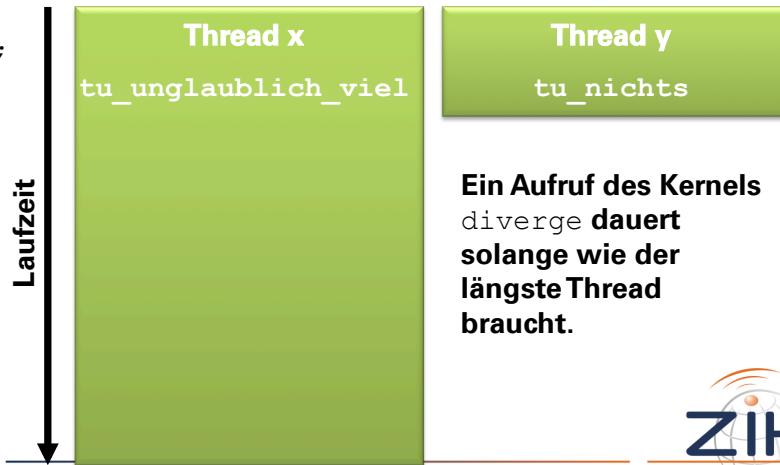
Thread 1 <code>*elem = 0.0f;</code>	Thread 2 <code>*elem = 0.0f;</code>	Thread 3 <code>*elem = 0.0f;</code>	Thread 4 <code>*elem = 0.0f;</code>
Thread 5 <code>*elem = 0.0f;</code>	Thread 6 <code>*elem = 0.0f;</code>	Thread ... <code>*elem = 0.0f;</code>	Thread ∞ <code>*elem = 0.0f;</code>

Exkurs: Thread-Divergenz



Werden verschiedene Teile des Kernel-Codes je nach Datenlage ausgeführt, so spricht man von Thread-Divergenz, da nun unterschiedliche Threads mit unterschiedlichen Daten unterschiedliche Dinge tun

```
__global__ void diverge( void *data ) {  
    if ( data[mythread] > zufallszahl )  
        tu_unglaublich_viel;  
    else  
        tu_nichts;  
}
```



Exkurs: Kernel-Aufrufe



Ab Compute Capability 2.x gibt es die Möglichkeit, mehrere Kernel parallel auszuführen. Diese werden in so genannten *Streams* organisiert.

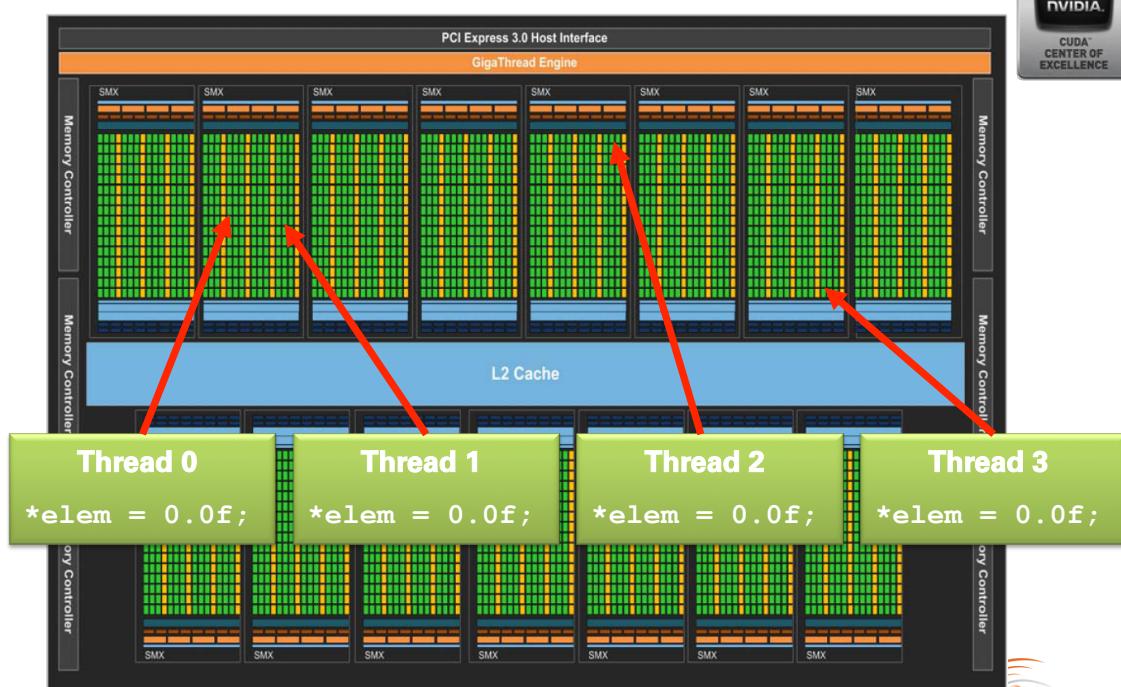
Diese Möglichkeit der Parallelisierung wird hier nicht besprochen, soll aber erwähnt werden.

Ab CUDA 5.0 und Kepler-Architektur gibt es zusätzlich die Möglichkeit, dass Kernel selbst mehrere Kernel ausführen. Sogar Rekursion ist möglich!

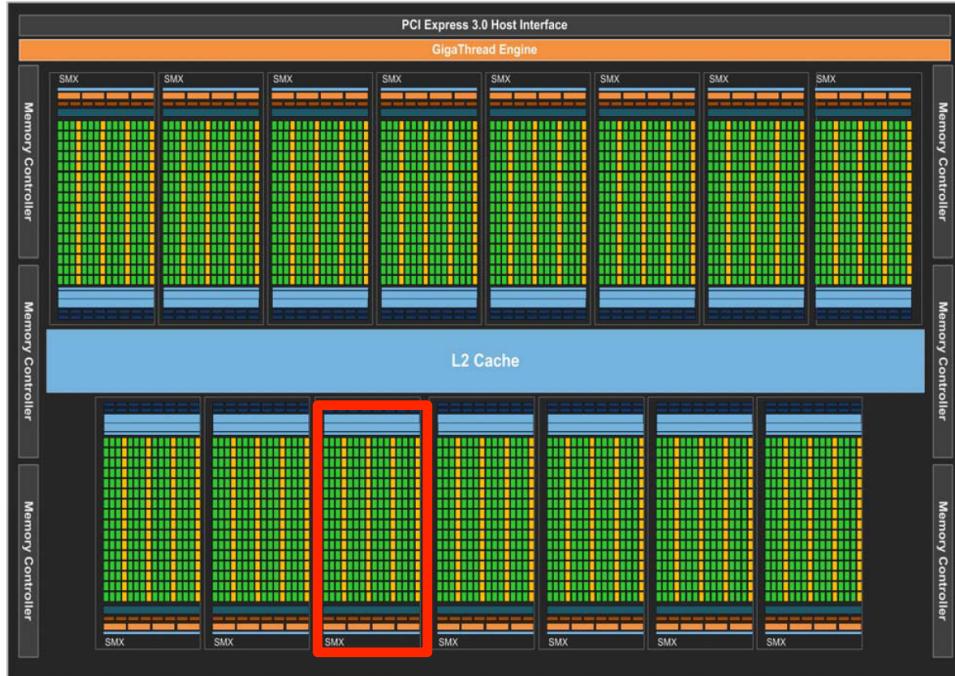
```
__global__ childkernel( void *data ) {...}  
  
__global__ parentkernel( void *data ) {  
    childkernel<<<...>>>( data ); //Dynamic Parallelism! YEAH!  
}  
  
...  
parentkernel<<<...>>>( data );
```

EIN AUSFLUG IN DIE ARCHITEKTUR VON GPUS

Kepler GK110 Architektur



Kepler GK110 Architektur



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 61



Streaming Multiprocessor



Threads, die parallel auf einem SMX ausgeführt werden sollen, sollten in so genannten Blöcken ausgeführt werden.

Sie können sich Speicher gemeinsam teilen und zusammen darauf zugreifen.

(hierzu später mehr)



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 62





WIE ERZEUGEN UND STEUERN WIR THREADS?



TECHNISCHE
UNIVERSITÄT
DRESDEN

Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 63



THREAD-BLOCKS



TECHNISCHE
UNIVERSITÄT
DRESDEN

Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 64



Thread-Blocks (Blöcke)



Ein Thread-Block ist eine Einheit von Threads, die gemeinsam auf einem Streaming Multiprocessor (SMX) parallel ausgeführt werden.

Die Reihenfolge, in der die Threads ausgeführt werden, nicht festgelegt, d.h. ein beliebiger Thread auf einem Core kann vor oder nach einem Thread auf einem anderen Core gestartet bzw. beendet werden.

Threads in einem Block können miteinander über gemeinsamen, schnellen Speicher (shared memory) Daten austauschen.

```
int nx; int ny; int nz; ...
dim3 block(nx, ny, nz); //nx, ny, nz = Größe des Blocks in 3D
kernel<<<1,block>>>(...); //Ruft nx*ny*nz threads auf, die in
                                //einem gemeinsamen Block laufen
```

Die Struktur dim3



```
struct dim3
{
    unsigned int x, y, z;
};
```

Thread-Blocks und Anzahl der Threads



Die Anzahl der Threads in einem Thread-Block wird durch das Produkt der Dimensionen des Thread-Blocks bestimmt:

```
dim3 block512(16,8,4) //Anzahl Threads ist 16*8*4 = 512
```

Die Anzahl der Threads in einem Block ist begrenzt. Ebenso ist die Anzahl der Threads in den einzelnen Dimensionen begrenzt.

Compute Capability	1.0	1.1	1.2	1.3	2.x	3.0
Max. Blockgröße in x,y		512			1024	
Max. Blockgröße in z				64		
Max. Threads pro Block		512			1024	

Was kann meine Hardware?



Wie bereits erwähnt, erhält man mit

```
cudaError_t cudaGetDeviceProperties( struct cudaDeviceProp *prop,  
                                     int      device );
```

Informationen zu seiner GPU. Diese enthalten auch die Informationen zu Thread-Blöcken:

```
struct cudaDeviceProp {  
    ...  
    int maxThreadsPerBlock; //Maximale Zahl der Threads pro Block  
    int maxThreadsDim[3];   //Maximale Größe eines Blocks in x,y,z  
    ...  
}
```



GRIDS OF BLOCKS

Kann ich mehr als 1024 Threads ausführen?



Es ist möglich, deutlich mehr Threads auszuführen. Dazu werden mehrere Thread-Blöcke parallel gestartet.

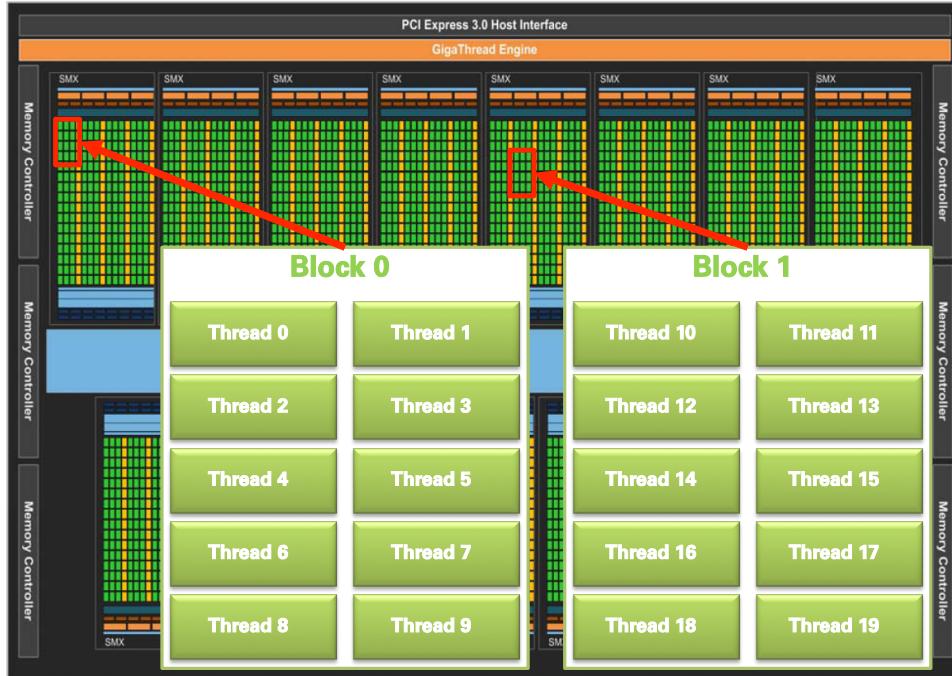
Hierzu erzeugt man Gitter (Grids) von Blöcken (Blocks).

Threads in unterschiedlichen Blöcken besitzen **keinen** gemeinsamen (shared) schnellen Speicher, über den sie Daten austauschen können.

Wie bei Blöcken ist die Abarbeitungsreihenfolge zufällig. Jeder Block im Gitter kann auf einem beliebigen (freien) Multiprozessor gestartet werden.

```
int nx; int ny; int nz; ...
dim3 grid(nx, ny, nz); //nx, ny, nz = Größe des Grids in 3D
kernel<<<grid,1>>>(...); //Ruft nx*ny*nz threads auf, jeden in
//einem eigenen, getrennt laufenden Block
```

Verteilung von Gittern



Gitter und Anzahl der Thread-Blöcke



Die Anzahl der Blöcke in einem Gitter wird durch das Produkt der Dimensionen des Gitters bestimmt:

```
dim3 grid512(16,8,4) //Anzahl Blöcke ist 16*8*4 = 512, CC ab 2.x
```

Die Anzahl der Blöcke in einem Gitter ist begrenzt. Ebenso ist die Anzahl der Blöcke in den einzelnen Dimensionen begrenzt.

Compute Capability	1.0	1.1	1.2	1.3	2.x	3.0
Dimension des Grids	2D (grid(nx,ny,1))				3D (grid(nx,ny,nz))	
Max. Gridgröße in x,y,z			65535			$2^{31}-1$
Max. Gridgröße			65535			$2^{31}-1$

Was kann meine Hardware?



Wiederum erhält man mit

```
cudaError_t cudaGetDeviceProperties( struct cudaDeviceProp *prop,  
                                     int      device );
```

Informationen zu seiner GPU. Diese enthalten Daten zu Gittern:

```
struct cudaDeviceProp {  
    ...  
    int maxGridSize[3]; //Maximale Größe eines Gitters in x,y,z  
    ...  
}
```



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 73



Zusammenfassung Gitter und Blöcke



Bei einem Kernelaufruf

```
dim3 block(bnx,bny,bnz);  
dim3 grid(gnx,gny,gnz);  
kernel<<<grid,block>>>(...);
```

werden $bnx \times bny \times bnz$ Threads in $gnx \times gny \times gnz$ Blöcken gestartet, d.h. es werden insgesamt $bnx \times bny \times bnz \times gnx \times gny \times gnz$ Threads aufgerufen (aber oft nicht gleichzeitig ausgeführt!).

Man kann einen Kernelaufruf daher auch verstehen als

```
kernel<<<Anzahl unabhängiger Thread-Blöcke,Threads pro Block>>>
```



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 74





THREADS IM BLOCK STARTEN

Threads in einem Block aufrufen

```
dim3 block1(8, 1, 1);  
kernel<<<1,block1>>>(...); //8 Threads in einem Block, alle in x  
oder auch einfach:  
kernel<<<1,8>>>(...); //8 Threads in einem Block
```



Threads in einem Block werden linear durchnummieriert.

Jeder Thread erhält eine eindeutige ID:

```
dim3 threadIdx
```

Ein Beispiel: Paralleles settozero im Block



```
#define NELEMENTS 512 //maximale Anzahl Threads pro Block
__global__ void settozero( float *elem ) {

    int tid = threadIdx.x; //Wir haben einen Vektor, d.h. eine
                           //eindimensionale Liste!
    elem[tid] = 0.0f;      //Jeder Thread setzt ein Element auf 0
}
```

Achtung! In diesem Beispiel nehmen wir an, dass unser Thread-Block mit

```
dim3 blockId(NELEMENTS,1,1)
settozero<<<1,blockId>>>(...)

aufgerufen wird. Dies ist equivalent zum Aufruf
settozero<<<1,NELEMENTS>>> (...)
```



Ein Beispiel: Paralleles settozero im Block



Um den gesamten Vektor parallel bearbeiten zu können, müssen wir Speicherplatz für den gesamten Vektor im GPU-Speicher reservieren

```
#define NELEMENTS 512 //maximale Anzahl Threads pro Block
int main( int argc, char *argv[] ) {

    float veca[NELEMENTS]; //Vektor im Host-Speicher
    float *devicemem;      //Pointer auf Vektor im GPU-Speicher

    //Jetzt reserviere Speicherplatz für den gesamten Vektor
    //auf der GPU, d.h. für NELEMENTS float-Variablen.
    cudaMalloc( (void**)&devicemem, NELEMENTS*sizeof(float) );
    ...

    cudaFree(devicemem);
    ...
}
```



Ein Beispiel: Paralleles settozero im Block



Nun müssen wir den Thread-Block definieren und aufrufen

```
#define NELEMENTS 512 //maximale Anzahl Threads pro Block
int main( int argc, char *argv[] ) {

    ...
    dim3 block1d(NELEMENTS,1,1); //1D-Block aus NELEMENTS Threads
    ...
    settozero<<<1,block1d>>>( devicemem ); //übergibt ein Array
    ...
}
```

Paralleles settozero im Block im Ganzen



```
#define NELEMENTS 512
__global__ void settozero( float *elem ) {
    int tid = threadIdx.x;
    elem[tid] = 0.0f;
}

int main( int argc, char *argv[] ) {

    float veca[NELEMENTS];
    float *devicemem;
    dim3 block1d(NELEMENTS,1,1);

    cudaMalloc( (void**)&devicemem, NELEMENTS*sizeof(float) );
    settozero<<<1,block1d>>>( devicemem );
    cudaMemcpy( veca, devicemem, NELEMENTS*sizeof(float), cudaMemcpyDeviceToHost );
    cudaFree( devicemem );

    return 0;
}
```

Threads in einem Block aufrufen



```
dim3 block1(8, 1, 1);  
kernel<<<1,block1>>>(...); //8 Threads in einem Block, alle in x  
oder auch einfach:  
kernel<<<1,8>>>(...); //8 Threads in einem Block
```

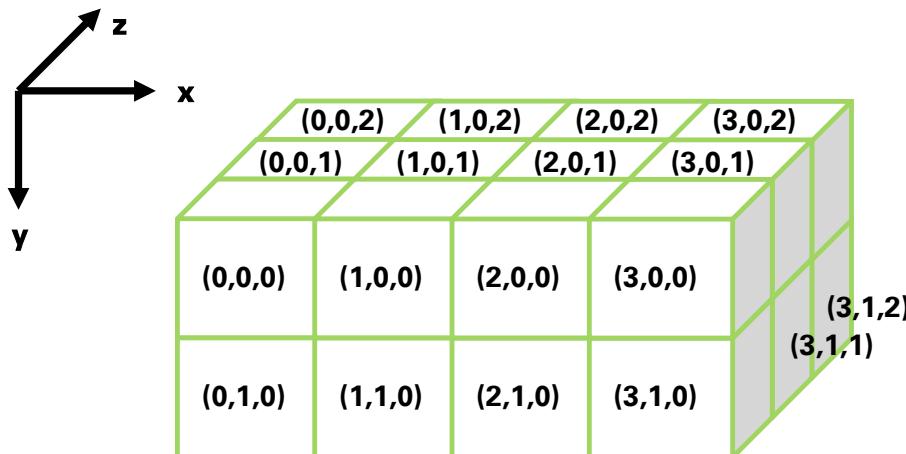
Andere Varianten:

```
dim3 block2(1, 8, 1);  
kernel<<<1,block2>>>(...); //8 Threads in einem Block, alle in y  
dim3 block3(1, 1, 8);  
kernel<<<1,block3>>>(...); //8 Threads in einem Block, alle in z  
dim3 block4(2, 1, 4);  
kernel<<<1,block4>>>(...); //8 Threads, 2 in x, 1 in y, 4 in z  
dim3 block5(2, 2, 2);  
kernel<<<1,block5>>>(...); //8 Threads, je 2 in x, y und z
```

Verteilung der Threads über den Block



```
dim3 block(4,2,3);  
kernel<<<1,block>>>(...);
```



New Threads on the Block

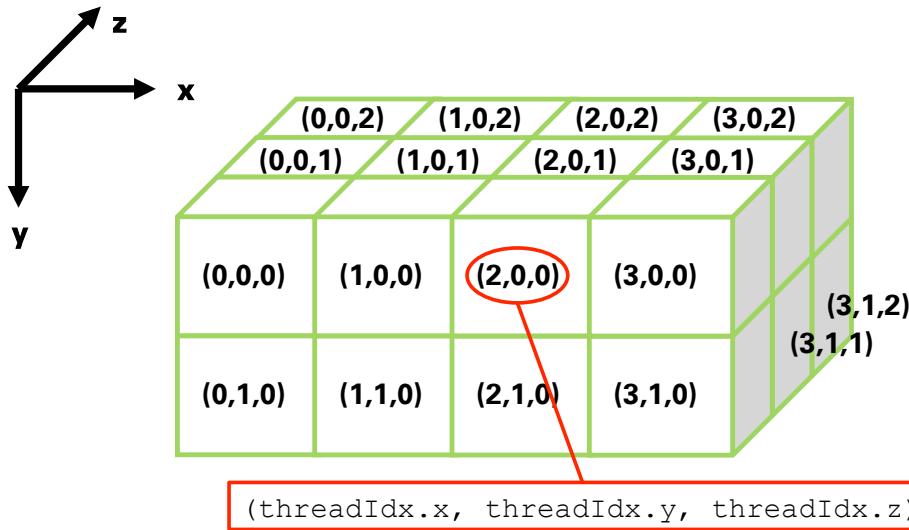


```
__global__ void kernel( void *data ) {  
    int tidx = threadIdx.x; //Position des Threads im Block in x  
    int tidy = threadIdx.y; //Position des Threads im Block in y  
    int tidz = threadIdx.z; //Position des Threads im Block in z  
}  
  
dim3 block(4,2,3);  
kernel<<<1,block>>>(data);  
ruft dann 24 = 4*2*3 Kernel auf mit  
(threadIdx.x, threadIdx.y, threadIdx.z):  
(0,0,0), (1,0,0), (2,0,0), (3,0,0), (0,1,0), (1,1,0), (2,1,0), (3,1,0),  
(0,0,1), (1,0,1), (2,0,1), (3,0,1), (0,1,1), (1,1,1), (2,1,1), (3,1,1),  
(0,0,2), (1,0,2), (2,0,2), (3,0,2), (0,1,2), (1,1,2), (2,1,2), (3,1,2)
```

Verteilung der Threads über den Block



```
dim3 block(4,2,3);  
kernel<<<1,block>>>(...);
```



Noch ein Beispiel: Matrix auf Null



```
#define DIMMAT 16
__global__ void settozero( float *elem ) {
    int tidx      = threadIdx.x;
    int tidy      = threadIdx.y;
    elem[tidx][tidy] = 0.0f;
}
int main( int argc, char *argv[] ) {

    float matrixa[DIMMAT][DIMMAT];
    float *devicemem;
    dim3 block2d(DIMMAT,DIMMAT,1);

    cudaMalloc( (void**)&devicemem, DIMMAT*DIMMAT*sizeof(float) );
    settozero<<<1,block2d>>>( devicemem );
    cudaMemcpy( veca, devicemem, DIMMAT*DIMMAT*sizeof(float), cudaMemcpyDeviceToHost );
    cudaFree( devicemem );
    return 0;
}
```



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 85



Alles auf Null — Linearisieren der Thread-ID



```
#define DIMX 4
#define DIMY 2
#define DIMZ 3
__global__ void settozero( float *elem ) {
    int tid = threadIdx.x + threadIdx.y * DIMX + threadIdx.z * DIMX * DIMY;
    elem[tid] = 0.0f;
}
int main( int argc, char *argv[] ) {

    float 3dstuff[DIMX*DIMY*DIMZ];
    float *devicemem;
    dim3 block3d(DIMX,DIMY,DIMZ);

    cudaMalloc( (void**)&devicemem, DIMX*DIMY*DIMZ*sizeof(float) );
    settozero<<<1,block3d>>>( devicemem );
    cudaMemcpy( veca, devicemem, DIMX*DIMY*DIMZ*sizeof(float), cudaMemcpyDeviceToHost );
    cudaFree( devicemem );
    return 0;
}
```



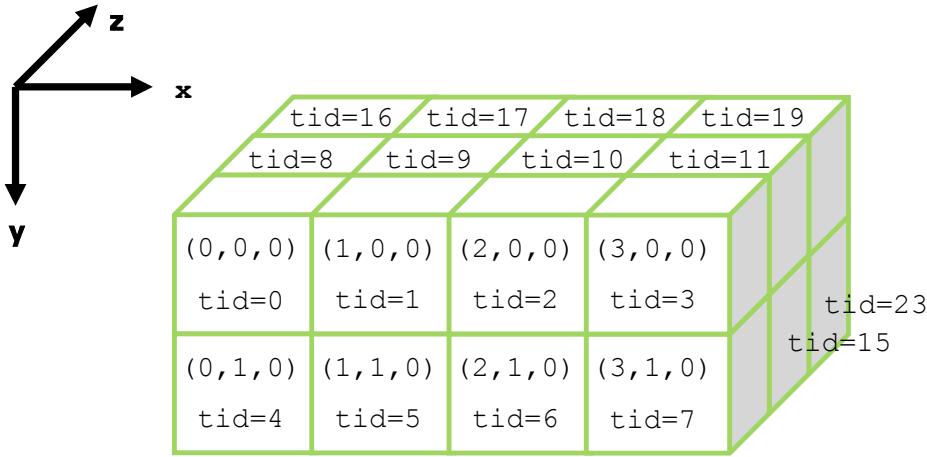
Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 86



Thread-IDs im Block



```
dim3 block3d(4,2,3);  
setzero<<<1,block3d>>>(...);
```

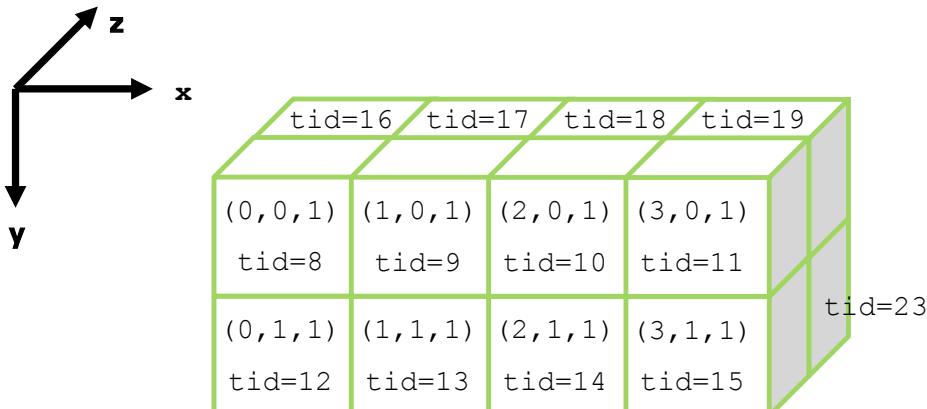


```
int tid = threadIdx.x + threadIdx.y * DIMX + threadIdx.z * DIMX * DIMY;
```

Thread-IDs im Block



```
dim3 block3d(4,2,3);  
setzero<<<1,block3d>>>(...);
```

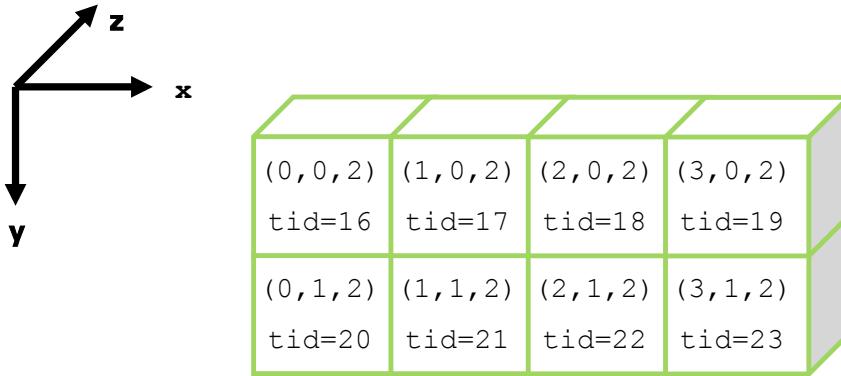


```
int tid = threadIdx.x + threadIdx.y * DIMX + threadIdx.z * DIMX * DIMY;
```

Thread-IDs im Block



```
dim3 block3d(4,2,3);  
setzero<<<1,block3d>>>(...);
```



```
int tid = threadIdx.x + threadIdx.y * DIMX + threadIdx.z * DIMX * DIMY;
```



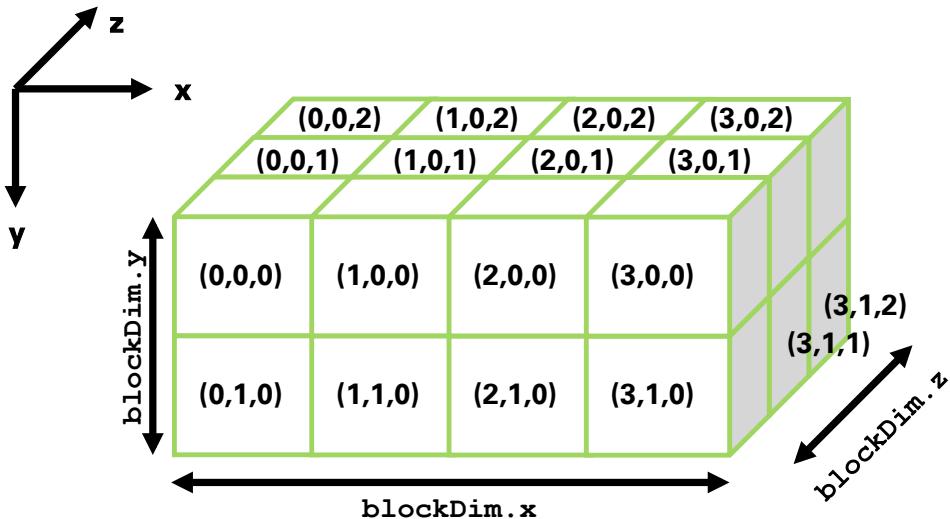
Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 89



CUDA hilft mit blockDim



```
dim3 block(4,2,3);  
kernel<<<1,block>>>(...);
```



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 90



CUDA hilft mit blockDim



```
#define DIMX 4
#define DIMY 2
#define DIMZ 3

__global__ void setzero( float *elem ) {
    int tid = threadIdx.x +
              threadIdx.y * blockDim.x +
              threadIdx.z * blockDim.x * blockDim.y;

    elem[tid] = 0.0f;
}

int main( int argc, char *argv[] ) {
...
dim3 block3d(DIMX,DIMY,DIMZ);
...
}
```



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 91



1-THREAD-BLÖCKE IM GITTER STARTEN



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 92



Threads (1-Thread-Blöcke) im Gitter starten



```
dim3 grid1(8, 1, 1);  
kernel<<<grid1,1>>>(...); //8 1-Thread-Blöcke, alle in x  
oder auch einfach:  
kernel<<<8,1>>>(...); //8 1-Thread-Blöcke
```

Ein Beispiel: Paralleles settozero im Gitter



```
#define NELEMENTS 65535 //maximale Anzahl Blöcke im Gitter  
__global__ void settozero( float *elem ) {  
  
    int tid = blockIdx.x; //Warum steht hier nicht threadIdx.x?  
    elem[tid] = 0.0f; //Jeder Thread setzt ein Element auf 0  
}
```

Achtung! In diesem Beispiel nehmen wir an, dass unser Thread-Block mit

```
dim3 grid1d(NELEMENTS,1,1)  
settozero<<<grid1d,1>>>(...)  
aufgerufen wird. Dies ist equivalent zum Aufruf  
settozero<<<NELEMENTS,1>>>(...)
```

Wir starten damit 65535 1-Thread-Blöcke! Jeder Thread hat in seinem Block ID 0!

Threads (1-Thread-Blöcke) im Gitter starten



```
dim3 grid1(8, 1, 1);  
kernel<<<grid1,1>>>(...); //8 1-Thread-Blöcke, alle in x  
oder auch einfach:  
kernel<<<8,1>>>(...); //8 1-Thread-Blöcke
```

Andere Varianten:

```
dim3 grid2(1, 8, 1);  
kernel<<<grid2,1>>>(...); //8 1-Thread-Blöcke, alle in y  
dim3 grid3(1, 1, 8);  
kernel<<<grid3,1>>>(...); //8 1-Thread-Blöcke, alle in z  
dim3 grid4(2, 1, 4);  
kernel<<<grid4,1>>>(...); //2 1-Thread-Blöcke in x, 1 in y, 4 in z  
dim3 grid5(2, 2, 2);  
kernel<<<grid5,1>>>(...); //Je 2 1-Thread-Blöcke in x, y und z
```

Block-ID und Gittergröße



Wie bei threadIdx gibt es auch eine ID für Blocks:

```
blockIdx.x //Block-ID im Gitter in x  
blockIdx.y //Block-ID im Gitter in y  
blockIdx.z //Block-ID im Gitter in z
```

Analog zur Blockgröße blockDim gibt es auch eine Gittergröße

```
gridDim.x //Größe des Gitters in x  
gridDim.y //Größe des Gitters in y  
gridDim.z //Größe des Gitters in z
```

Wiederholung: Matrix auf Null mit Block-Gitter



```
#define DIMMAT 255
__global__ void settozero( float *elem ) {
    int tidx      = blockIdx.x;
    int tidy      = blockIdx.y;
    elem[tidx][tidy] = 0.0f;
}
int main( int argc, char *argv[] ) {

    float matrixa[DIMMAT][DIMMAT];
    float *devicemem;
    dim3 grid2d(DIMMAT,DIMMAT,1);

    cudaMalloc( (void**)&devicemem, DIMMAT*DIMMAT*sizeof(float) );
    settozero<<<grid2d,1>>>( devicemem ); //startet DIMMAT*DIMMAT 1-Thread-Blöcke
    cudaMemcpy( veca, devicemem, DIMMAT*DIMMAT*sizeof(float), cudaMemcpyDeviceToHost );
    cudaFree( devicemem );
    return 0;
}
```



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 97



Threads in Blöcken vs. 1-Thread-Blöcke in Gittern

```
#define DIMX 4
#define DIMY 2
#define DIMZ 3
__global__ void settozero( float *elem ) {
    int tid = threadIdx.x + 
              threadIdx.y * blockDim.x +
              threadIdx.z * blockDim.x * blockDim.y;

    elem[tid] = 0.0f;
}

int main( int argc, char *argv[] ) {
...
dim3 block3d(DIMX,DIMY,DIMZ);
...
settozero<<<1,block3d>>>( devicemem )
...
}
```



```
int tid = blockIdx.x + 
          blockIdx.y * gridDim.x +
          blockIdx.z * gridDim.x * gridDim.y;

dim3 grid3d(DIMX,DIMY,DIMZ);
settozero<<<grid3d,1>>>( devicemem )
```



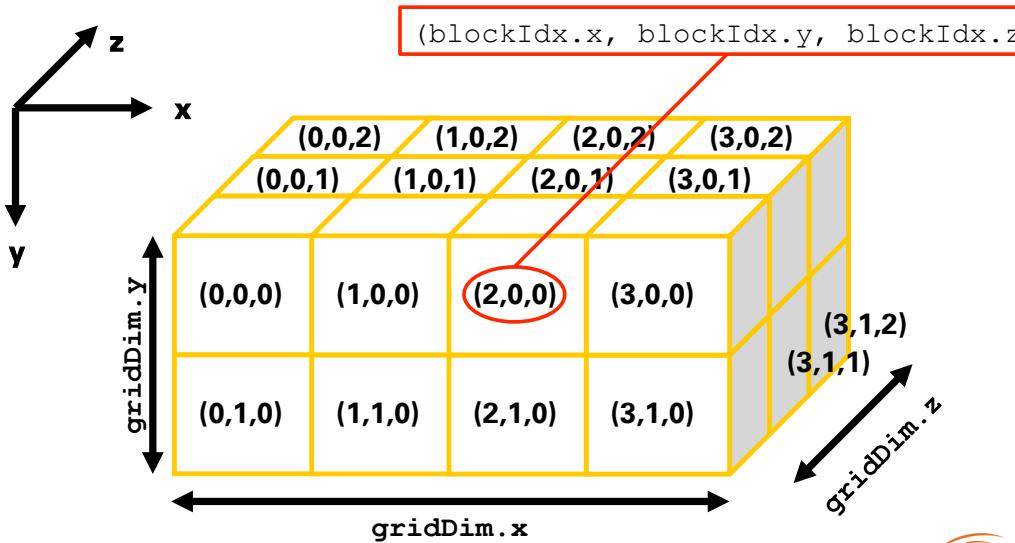
Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 98



Blockverteilung im Gitter



```
dim3 grid(4,2,3);  
kernel<<<grid, 1>>> (...);
```



BLÖCKE IM GITTER STARTEN

Noch einmal: Paralleles settozero



Das Problem mit unserem parallelen settozero ist, dass wir im Block auf maxThreadsPerBlock (512 oder 1024) Threads beschränkt sind.

Das heißt, unser Code kann nur Vektoren (Arrays) mit maximal so vielen Elementen wie Threads in einen Block passen auf Null setzen.

Mit Gittern können wir immerhin schon `maxGridSize[1]` 1-Thread-Blöcke starten, also insgesamt 65535 oder $2^{31}-1$ Threads., je nach Compute Capability.

Jetzt wollen wir Gitter aus Blöcken bauen, bei denen jeder Block mehrere Threads enthält. So können wir zunächst einmal bis zu

$$\text{maxThreadsPerBlock} * \text{maxGridSize}[1] \leq 1024 * (2^{31}-1)$$

Kernelaufufe starten.

Paralleles settozero im 1D Blockgitter



```
#define NTHREADS 512    //maximale Anzahl Threads pro Block
#define NBLOCKS  65535 //maximale Anzahl von Blöcken im Gitter
__global__ void settozero( float *elem ) {

    int tid    = threadIdx.x + blockIdx.x * blockDim.x;
    elem[tid] = 0.0f;
}

dim3 block1d(NTHREADS,1,1);
dim3 grid1d(NBLOCKS,1,1);
settozero<<<grid1d,block1d>>>(...)
```

Dies ist equivalent zum Aufruf

`settozero<<<NBLOCKS, NTHREADS>>>(...)`

Paralleles settozero im 1D Blockgitter



Gitter	Blöcke					
blockIdx.x	threadIdx.x	threadIdx.x	threadIdx.x		threadIdx.x	
0	0	1	2	...	511	
1	0	1	2	...	511	
2	0	1	2	...	511	
...	0	1	2	...	511	
65534	0	1	2	...	511	

Paralleles settozero im 1D Blockgitter



Gitter	Blöcke					
blockIdx.x	threadIdx.x	threadIdx.x	threadIdx.x		threadIdx.x	
0	0	1	2	...	511	
1	tid=0	tid=1	tid=2	...	tid=511	
2	tid=512	tid=513	tid=514	...	tid=1023	
3	tid=1024	tid=1025	tid=1026	...	tid=1535	
4	tid=1536	tid=1537	tid=1538	...	tid=2047	

$tid = threadIdx.x + blockIdx.x * blockDim.x;$

Paralleles settozero im 1D Blockgitter



```
int main( int argc, char *argv[] ) {  
  
    float veca[NTHREADS*NBLOCKS];  
    float *devicemem;  
    dim3 block1d(NTHREADS,1,1);  
    dim3 grid1d(NBLOCKS,1,1);  
  
    cudaMalloc( (void**)&devicemem, NTHREADS*NBLOCKS*sizeof(float) );  
    settozero<<<grid1d,block1d>>>( devicemem ); //übergibt ein Array  
    cudaMemcpy( veca, devicemem, NTHREADS*NBLOCKS*sizeof(float),  
               cudaMemcpyDeviceToHost ); //kopiert Array  
    cudaFree(devicemem);  
    return 0;  
}
```

Tabelle zu IDs und Größen



Typ	ID	Größe
Thread	threadIdx	1
Block	blockIdx	blockDim
Grid	-	gridDim

Blöcke im Gitter starten



```
dim3 grid1(8, 1, 1);  
dim3 block(nx,ny,nz);  
kernel<<<grid1,block>>>(...); //8 Blöcke, alle in x  
oder auch einfach:  
kernel<<<8,block>>>(...); //8 Blöcke
```

Blöcke im Gitter starten

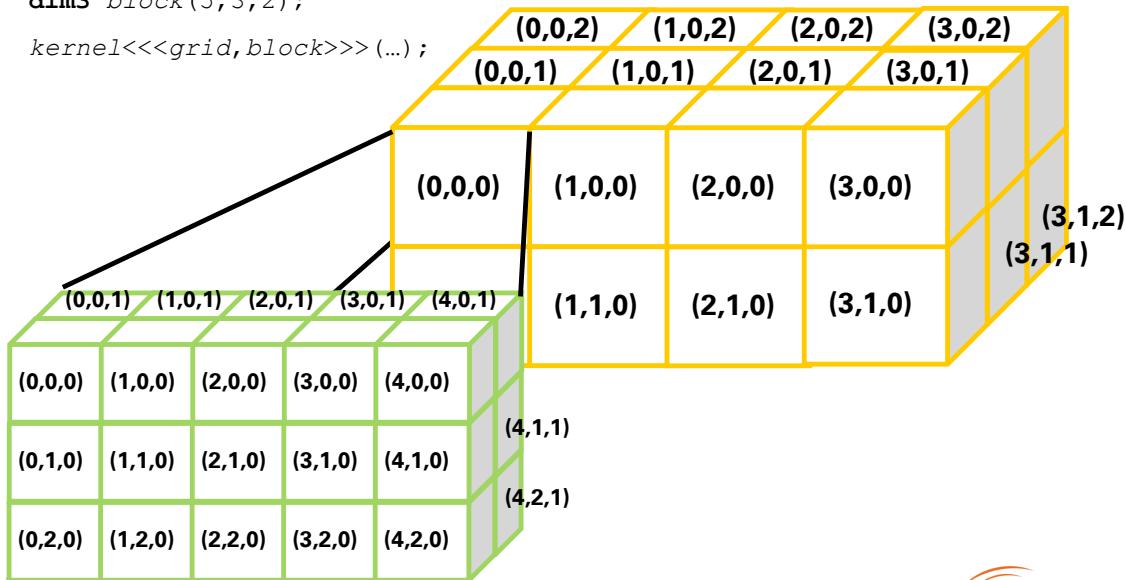


```
dim3 grid1(8, 1, 1);  
dim3 block(nx,ny,nz);  
kernel<<<grid1,block>>>(...); //8 Blöcke, alle in x  
oder auch einfach:  
kernel<<<8,block>>>(...); //8 Blöcke  
Andere Varianten:  
dim3 grid2(1, 8, 1);  
kernel<<<grid2,block>>>(...); //8 Blöcke, alle in y  
dim3 grid3(1, 1, 8);  
kernel<<<grid3,block>>>(...); //8 Blöcke, alle in z  
dim3 grid4(2, 1, 4);  
kernel<<<grid4,block>>>(...); //2 Blöcke in x, 1 in y, 4 in z  
dim3 grid5(2, 2, 2);  
kernel<<<grid5,block>>>(...); //Je 2 Blöcke in x, y und z
```

Threads in Blöcken in Gittern



```
dim3 grid(4,2,3);  
dim3 block(5,3,2);  
kernel<<<grid,block>>> (...);
```



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 109



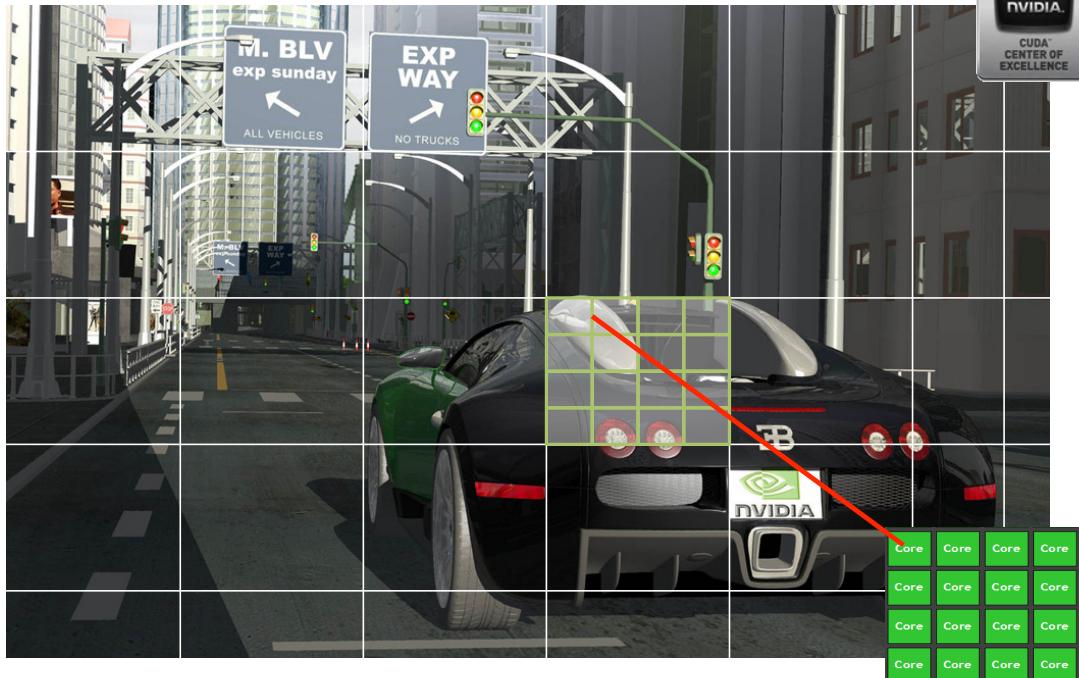
PARALLELE TECHNIKEN: DOMAIN DECOMPOSITION



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 110



Gebietszerlegung (Domain Decomposition)



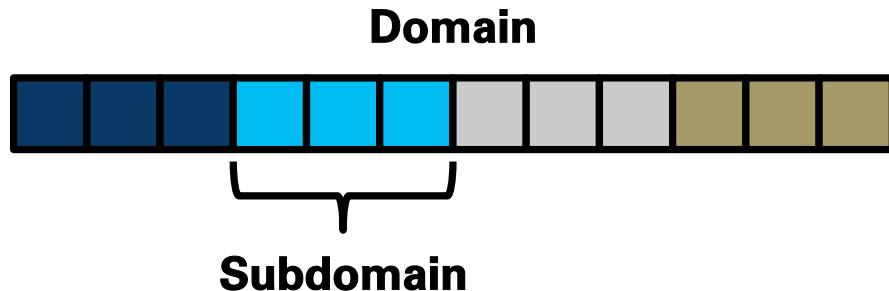
Domain Decomposition in 1D



Domain



Domain Decomposition in 1D



„Partitioning“

„Subdivision“

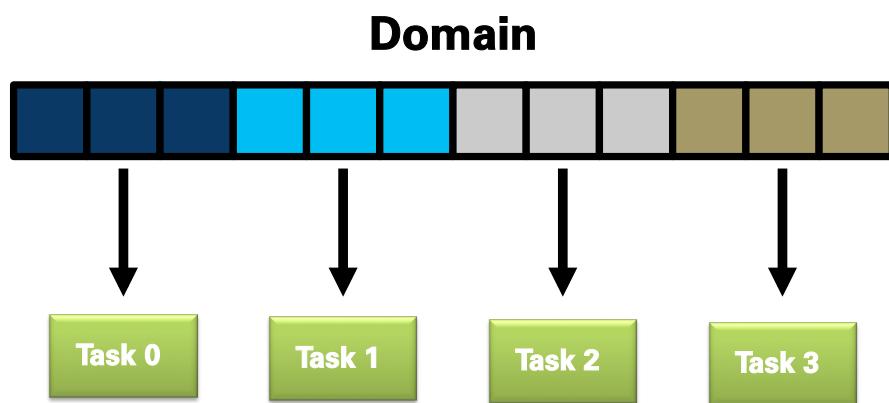


TECHNISCHE
UNIVERSITÄT
DRESDEN

Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 113



Domain Decomposition in 1D

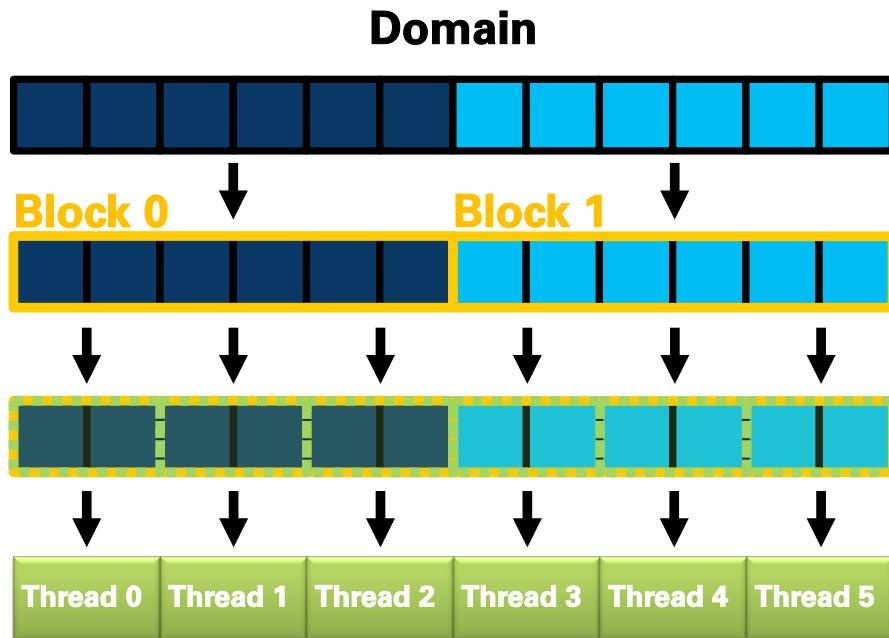


TECHNISCHE
UNIVERSITÄT
DRESDEN

Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 114



Domain Decomposition in 1D in CUDA



Domain Decomposition in 1D in CUDA



```
__global__ void setzero( float *elem, int NELEMENTS ) {  
  
    int subDim_x      = NELEMENTS / ( blockDim.x * gridDim.x ); //Anz. Elemente in Subarray  
    int tid           = threadIdx.x + blockIdx.x * blockDim.x; //Thread-ID  
    int subStartIdx = tid * subDim_x                         //Startindex im Array  
    int i;                                         //Laufindex im Array  
  
    for ( i=0; i++; i<subDimx ) elem[i+subStartIdx] = 0.0f;  
}  
  
int main( int argc, char *argv[] ) {  
  
    const int NELEMENTS = 128000;  
    float veca[NELEMENTS];  
    ...  
    setzero<<10,128>>( devicemem, NELEMENTS ); //subDim_x = 100  
    ...  
}
```

Die Welt in 3D



```
__global__ void WorkOnSubdomain( data *domain, dim3 domainDim ) {  
  
    dim3 subDim;      //Größe der Subdomain in x,y,z  
    dim3 tid;         //Thread-ID in x,y,z  
    dim3 subStartIdx; //Anfangsposition der Daten der Subdomain in der Domain  
  
    subDim.x          = domainDim.x / ( blockDim.x * gridDim.x ); //Achtung, Integerdivision  
    subDim.y          = domainDim.y / ( blockDim.y * gridDim.y ); //Achtung, Integerdivision  
    subDim.z          = domainDim.z / ( blockDim.z * gridDim.z ); //Achtung, Integerdivision  
    tid.x            = threadIdx.x + blockIdx.x * blockDim.x;  
    tid.y            = threadIdx.y + blockIdx.y * blockDim.y;  
    tid.z            = threadIdx.z + blockIdx.z * blockDim.z;  
  
    subStartIdx.x = tid.x * subDim.x;  
    subStartIdx.y = tid.y * subDim.y;  
    subStartIdx.z = tid.z * subDim.z;  
  
    //Schleifen in x,y,z oder ähnliches...  
}
```



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 117



Domain Decomposition in 1D in CUDA



domainDim.x = 12



gridDim.x = 2



blockDim.x = 3



subDim.x = 2

startIndex.x = 6



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 118



Die „CUDA by Example“-Variante



Domain



Die „CUDA by Example“-Variante



Domain



i

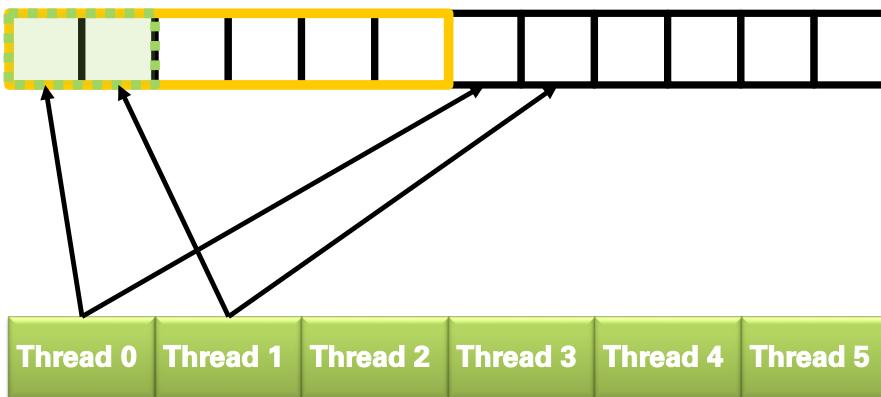
$i + blockDim.x * gridDim.x$



Die „CUDA by Example“-Variante



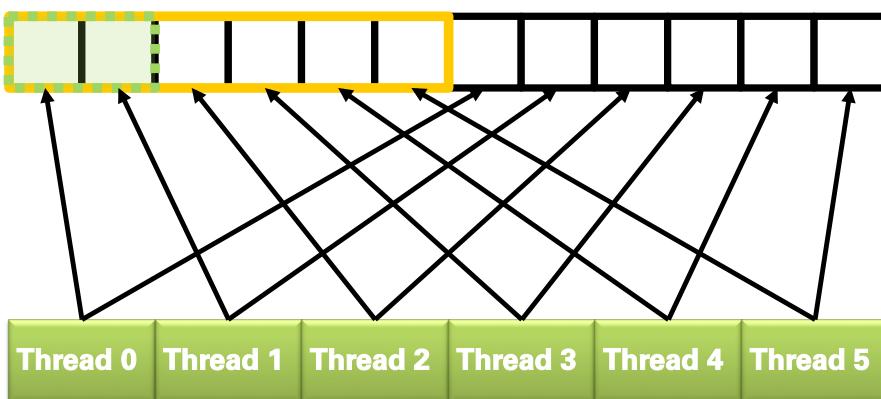
Domain



Die „CUDA by Example“-Variante



Domain



Domain Decomposition in 1D „CUDA by Example“



```
__global__ void setzero( float *elem, int NELEMENTS ) {  
  
    int tid      = threadIdx.x + blockIdx.x * blockDim.x; //Thread-ID  
    int idxinc  = blockDim.x * gridDim.x;                  //Index-Inkrement im Array  
    int i        = tid;                                  //Laufindex im Array  
  
    while ( i<NELEMENTS ) {  
        elem[i] = 0.0f;  
        i       += idxinc;  
    }  
}
```



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 123



BEISPIELE FÜR DOMAIN DECOMPOSITION

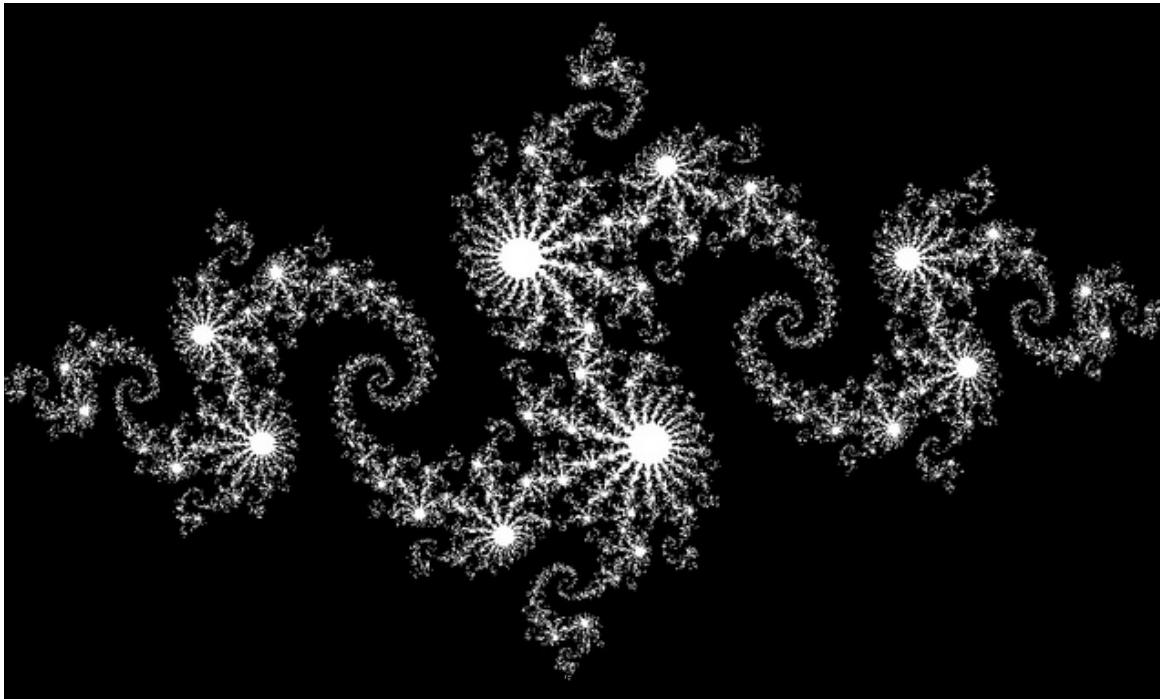


Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 124





Julia-Menge



Paralleles Suchen



Lorem ipsum

Lorem ipsum dolor sit amet, consectetur adipisici elit, sed eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquid ex ea commodo consequat. Quis aute iure reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint obcaecat cupiditat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

--

Duis autem vel eum iriure dolor in hendrerit in vulputate velit esse molestie consequat, vel illum dolore eu feugiat nulla facilisis at vero eros et accumsan et iusto odio dignissim qui blandit praesent luptatum zzril delenit augue duis dolore te feugait nulla facilisi. Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed diam nonummy nibh euismod tincidunt ut laoreet dolore magna aliquam erat volutpat.

Ut wisi enim ad minim veniam, quis nostrud exerci tation ullamcorper suscipit lobortis nisl ut aliquip ex ea commodo consequat. Duis autem vel eum iriure dolor in hendrerit in vulputate velit esse molestie consequat, vel illum dolore eu feugiat nulla facilisis at vero eros et accumsan et iusto odio dignissim qui blandit praesent luptatum zzril delenit augue duis dolore te feugait nulla facilisi.

Nam liber tempor cum soluta nobis eleifend option congue nihil imperdiet doming id quod mazim placerat facer possim assum. Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed diam nonummy nibh euismod tincidunt ut laoreet dolore magna aliquam erat volutpat. Ut wisi enim ad minim veniam, quis nostrud exerci tation ullamcorper suscipit lobortis nisl ut aliquip ex ea commodo consequat.

Duis autem vel eum iriure dolor in hendrerit in vulputate velit esse molestie consequat, vel illum dolore eu feugiat nulla facilisis.

Paralleles Suchen



Lore ipsum

Lore ipsum dolor sit amet, consectetur adipisici elit, sed eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquid ex ea commodo consequat. Quis aute iure reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint obcaecat cupiditat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

--

Duis autem vel eum iriure dolor in hendrerit in vulputate velit esse molestie consequat, vel illum dolore eu feugiat nulla facilisis at vero eros et accumsan et iusto odio dignissim qui blandit praesent luptatum zzril delenit augue duis dolore te feugait nulla facilisi. Lorem ipsum dolor sit amet, consectetur adipisciing elit, sed diam nonumy nibh euismod tincidunt ut laoreet dolore magna aliquam erat volutpat.

Ut wisi enim ad minim veniam, quis nostrud exerci tation ullamcorper suscipit lobortis nisl ut aliquip ex ea commodo consequat. Duis autem vel eum iriure dolor in hendrerit in vulputate velit esse molestie consequat, vel illum dolore eu feugiat nulla facilisis at vero eros et accumsan et iusto odio dignissim qui blandit praesent luptatum zzril delenit augue duis dolore te feugait nulla facilisi.

Nam liber tempor cum soluta nobis eleifend option congue nihil imperdiet doming id quod mazim placerat facer possim assum. Lorem ipsum dolor sit amet, consectetur adipisciing elit, sed diam nonumy nibh euismod tincidunt ut laoreet dolore magna aliquam erat volutpat. Ut wisi enim ad minim veniam, quis nostrud exerci tation ullamcorper suscipit lobortis nisl ut aliquip ex ea commodo consequat.

Duis autem vel eum iriure dolor in hendrerit in vulputate velit esse molestie consequat, vel illum dolore eu feugiat nulla facilisis.



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 127



Paralleles Suchen



Lore ipsum

Lore ipsum dolor sit amet, consectetur adipisici elit, sed eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquid ex ea commodo consequat. Quis aute iure reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint obcaecat cupiditat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

--

Duis autem vel eum iriure dolor in hendrerit in vulputate velit esse molestie consequat, vel illum dolore eu feugiat nulla facilisis at vero eros et accumsan et iusto odio dignissim qui blandit praesent luptatum zzril delenit augue duis dolore te feugait nulla facilisi. Lorem ipsum dolor sit amet, consectetur adipisciing elit, sed diam nonumy nibh euismod tincidunt ut laoreet dolore magna aliquam erat volutpat.

Ut wisi enim ad minim veniam, quis nostrud exerci tation ullamcorper suscipit lobortis nisl ut aliquip ex ea commodo consequat. Duis autem vel eum iriure dolor in hendrerit in vulputate velit esse molestie consequat, vel illum dolore eu feugiat nulla facilisis at vero eros et accumsan et iusto odio dignissim qui blandit praesent luptatum zzril delenit augue duis dolore te feugait nulla facilisi.

Nam liber tempor cum soluta nobis eleifend option congue nihil imperdiet doming id quod mazim placerat facer possim assum. Lorem ipsum dolor sit amet, consectetur adipisciing elit, sed diam nonumy nibh euismod tincidunt ut laoreet dolore magna aliquam erat volutpat. Ut wisi enim ad minim veniam, quis nostrud exerci tation ullamcorper suscipit lobortis nisl ut aliquip ex ea commodo consequat.

Duis autem vel eum iriure dolor in hendrerit in vulputate velit esse molestie consequat, vel illum dolore eu feugiat nulla facilisis.



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 128



Paralleles Suchen



Lorem ipsum

Lorem ipsum dolor sit amet, consectetur adipisici elit, sed eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquid ex ea commodo consequat. Quis aute iure reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint obcaecat cupiditat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

--

Duis autem vel eum iriure dolor in hendrerit in vulputate velit esse molestie consequat, vel illum dolore eu feugiat nulla facilisis at vero eros et accumsan et iusto odio dignissim qui blandit praesent luptatum zzril delenit augue duis dolore te feugait nulla facilisi. Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed diam nonummy nibh euismod tincidunt ut laoreet dolore magna aliquam erat volutpat.

Ut wisi enim ad minim veniam, quis nostrud exerci tation ullamcorper suscipit lobortis nisl ut aliquip ex ea commodo consequat. Duis autem vel eum iriure dolor in hendrerit in vulputate velit esse molestie consequat, vel illum dolore eu feugiat nulla facilisis at vero eros et accumsan et iusto odio dignissim qui blandit praesent luptatum zzril delenit augue duis dolore te feugait nulla facilisi.

Nam liber tempor cum soluta nobis eleifend option congue nihil imperdiet doming id quod mazim placerat facer possim assum. Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed diam nonummy nibh euismod tincidunt ut laoreet dolore magna aliquam erat volutpat. Ut wisi enim ad minim veniam, quis nostrud exerci tation ullamcorper suscipit lobortis nisl ut aliquip ex ea commodo consequat.

Duis autem vel eum iriure dolor in hendrerit in vulputate velit esse molestie consequat, vel illum dolore eu feugiat nulla facilisis.



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 129



Paralleles Suchen



Lorem ipsum

Lorem ipsum dolor sit amet, consectetur adipisici elit, sed eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquid ex ea commodo consequat. Quis aute iure reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint obcaecat cupiditat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

--

Duis autem vel eum iriure dolor in hendrerit in vulputate velit esse molestie consequat, vel illum dolore eu feugiat nulla facilisis at vero eros et accumsan et iusto odio dignissim qui blandit praesent luptatum zzril delenit augue duis dolore te feugait nulla facilisi. Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed diam nonummy nibh euismod tincidunt ut laoreet dolore magna aliquam erat volutpat.

Ut wisi enim ad minim veniam, quis nostrud exerci tation ullamcorper suscipit lobortis nisl ut aliquip ex ea commodo consequat. Duis autem vel eum iriure dolor in hendrerit in vulputate velit esse molestie consequat, vel illum dolore eu feugiat nulla facilisis at vero eros et accumsan et iusto odio dignissim qui blandit praesent luptatum zzril delenit augue duis dolore te feugait nulla facilisi.

Nam liber tempor cum soluta nobis eleifend option congue nihil imperdiet doming id quod mazim placerat facer possim assum. Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed diam nonummy nibh euismod tincidunt ut laoreet dolore magna aliquam erat volutpat. Ut wisi enim ad minim veniam, quis nostrud exerci tation ullamcorper suscipit lobortis nisl ut aliquip ex ea commodo consequat.

Duis autem vel eum iriure dolor in hendrerit in vulputate velit esse molestie consequat, vel illum dolore eu feugiat nulla facilisis.



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 130





Lorem ipsum dolor sit amet, consectetur adipisici elit, sed eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquid ex ea commodi consequat. Quis aute iure reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint obcaecat cupiditat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

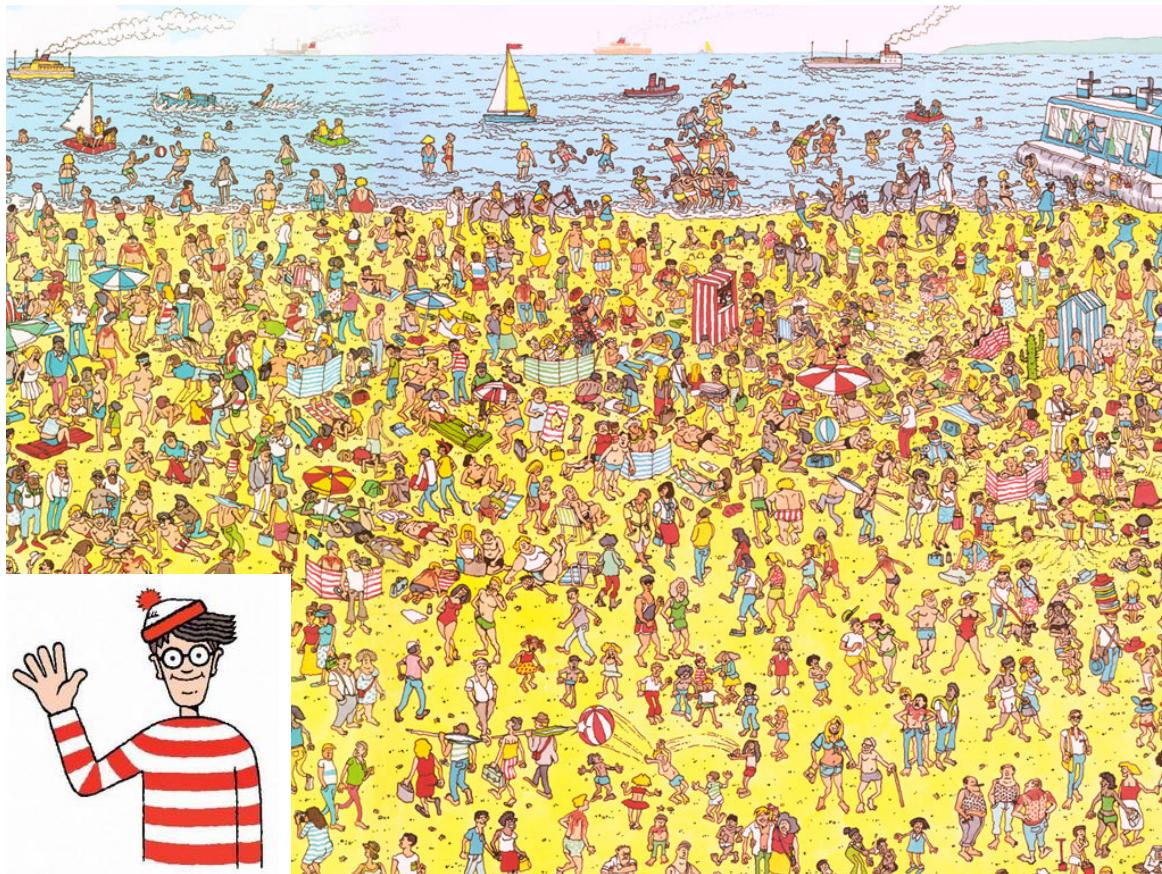
--

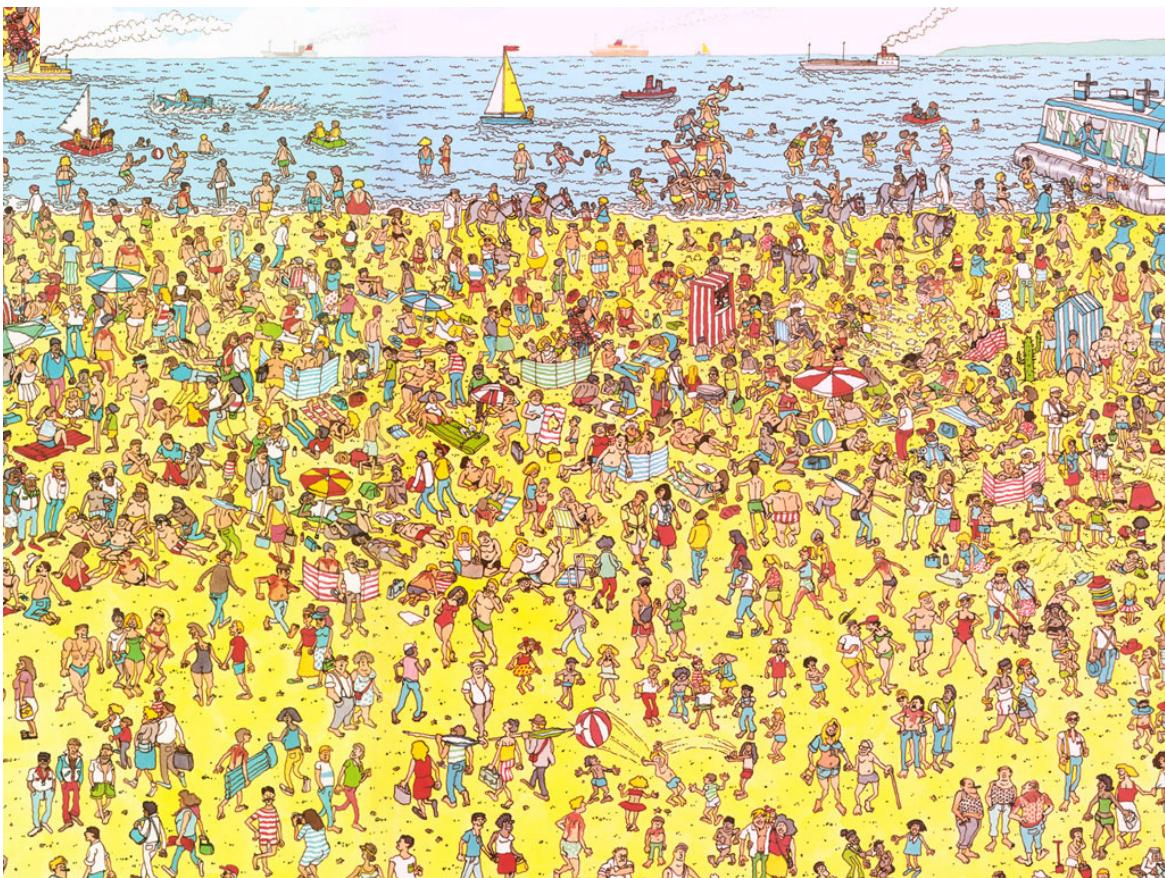
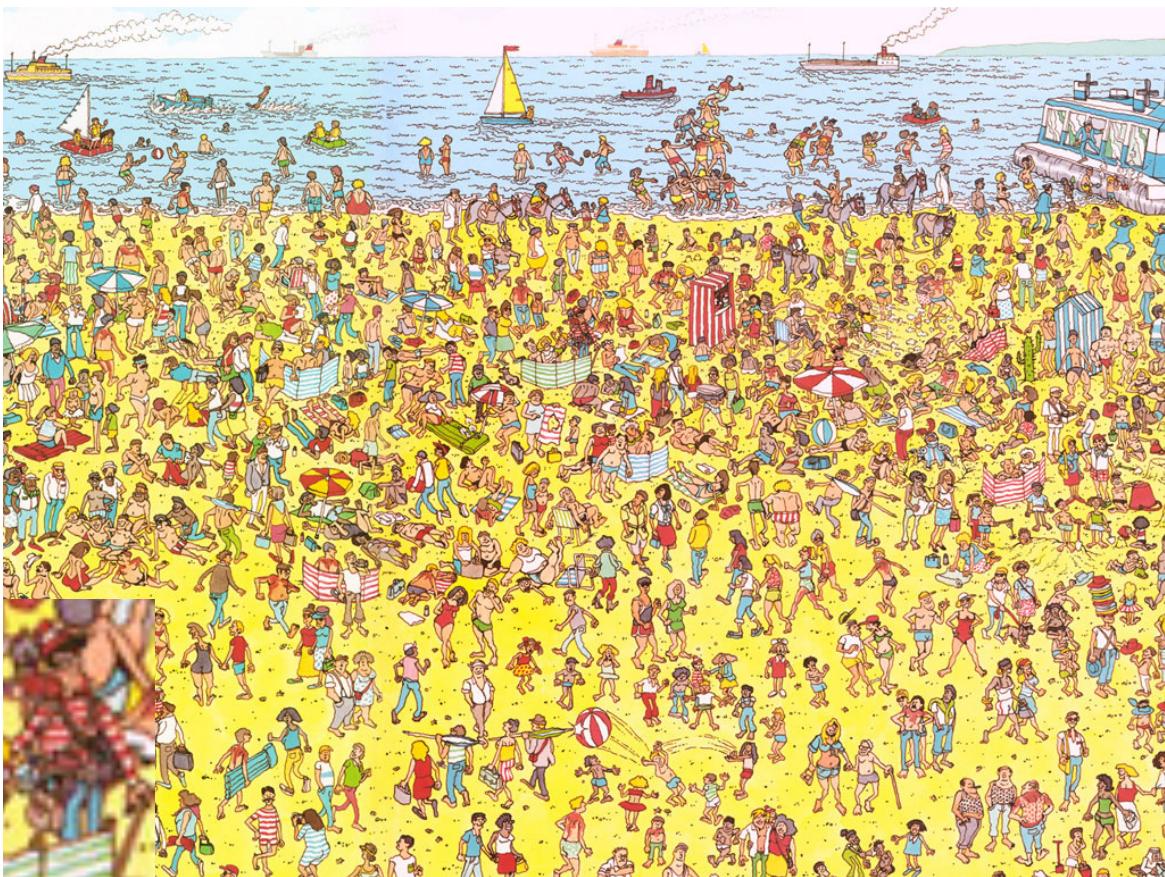
Duis autem vel eum iriure dolor in hendrerit in vulputate velit esse molestie consequat, vel illum dolore eu feugiat nulla facilisis at vero eros et accumsan et iusto odio dignissim qui blandit praesent luptatum zzril delenit augue duis dolore te feugait nulla facilisi. **Lorem ipsum** dolor sit amet, consectetur adipiscing elit, sed diam nonummy nibh euismod tincidunt ut laoreet dolore magna aliquam erat volutpat.

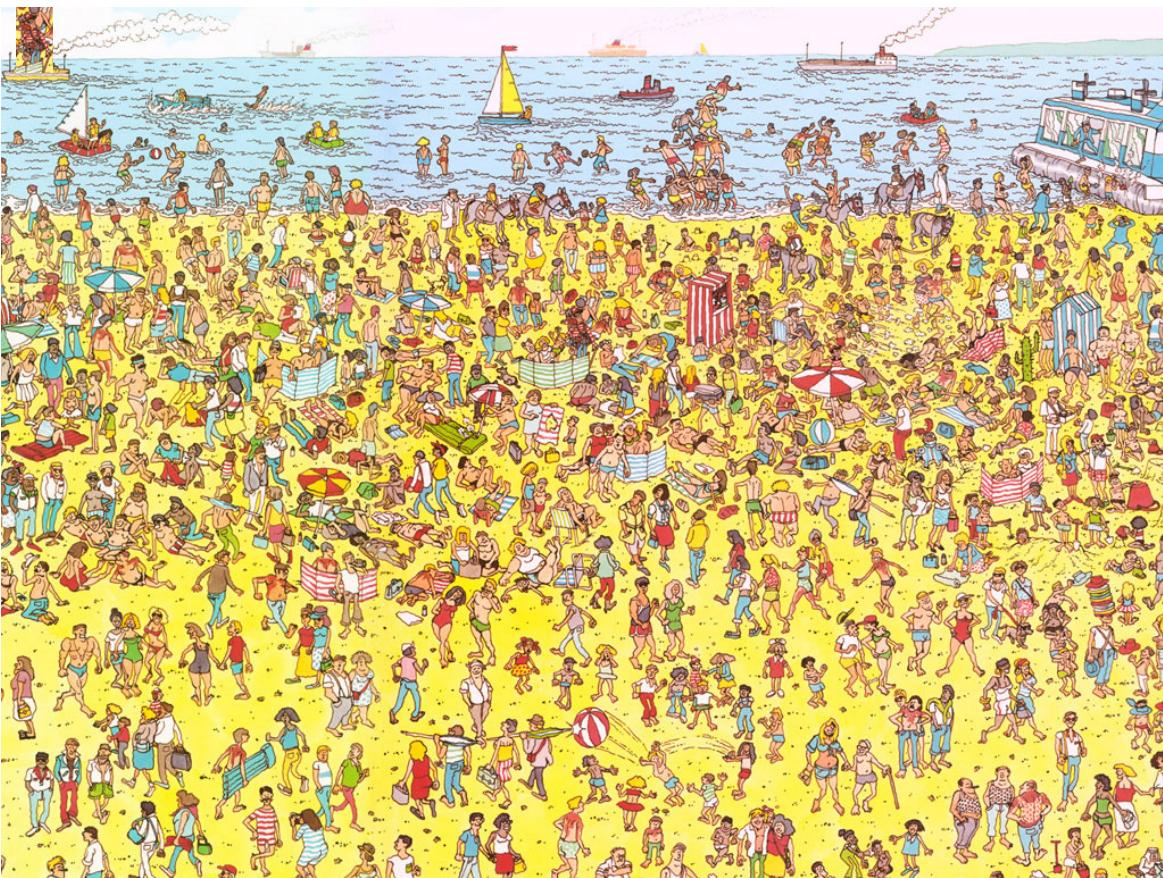
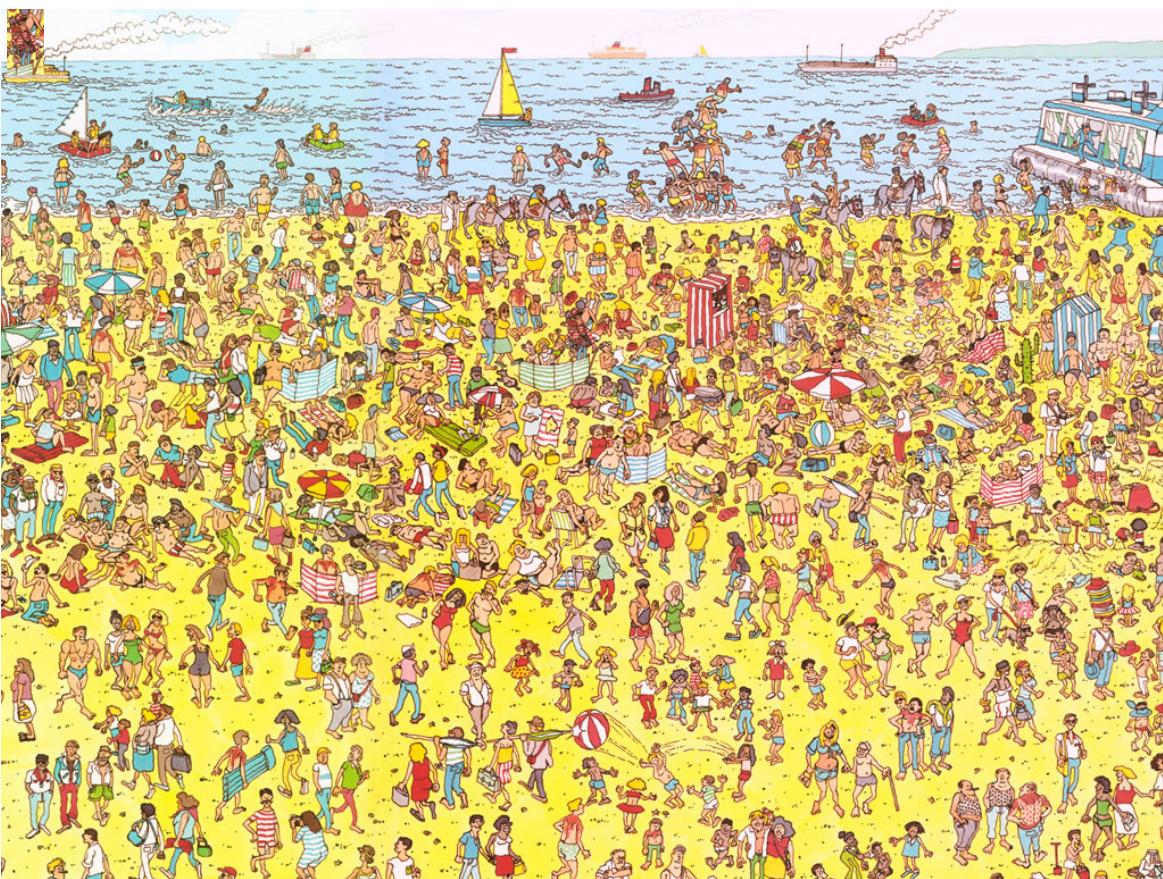
Ut wisi enim ad minim veniam, quis nostrud exerci tation ullamcorper suscipit lobortis nisl ut aliquip ex ea commodo consequat. Duis autem vel eum iriure dolor in hendrerit in vulputate velit esse molestie consequat, vel illum dolore eu feugiat nulla facilisis at vero eros et accumsan et iusto odio dignissim qui blandit praesent luptatum zzril delenit augue duis dolore te feugait nulla facilisi.

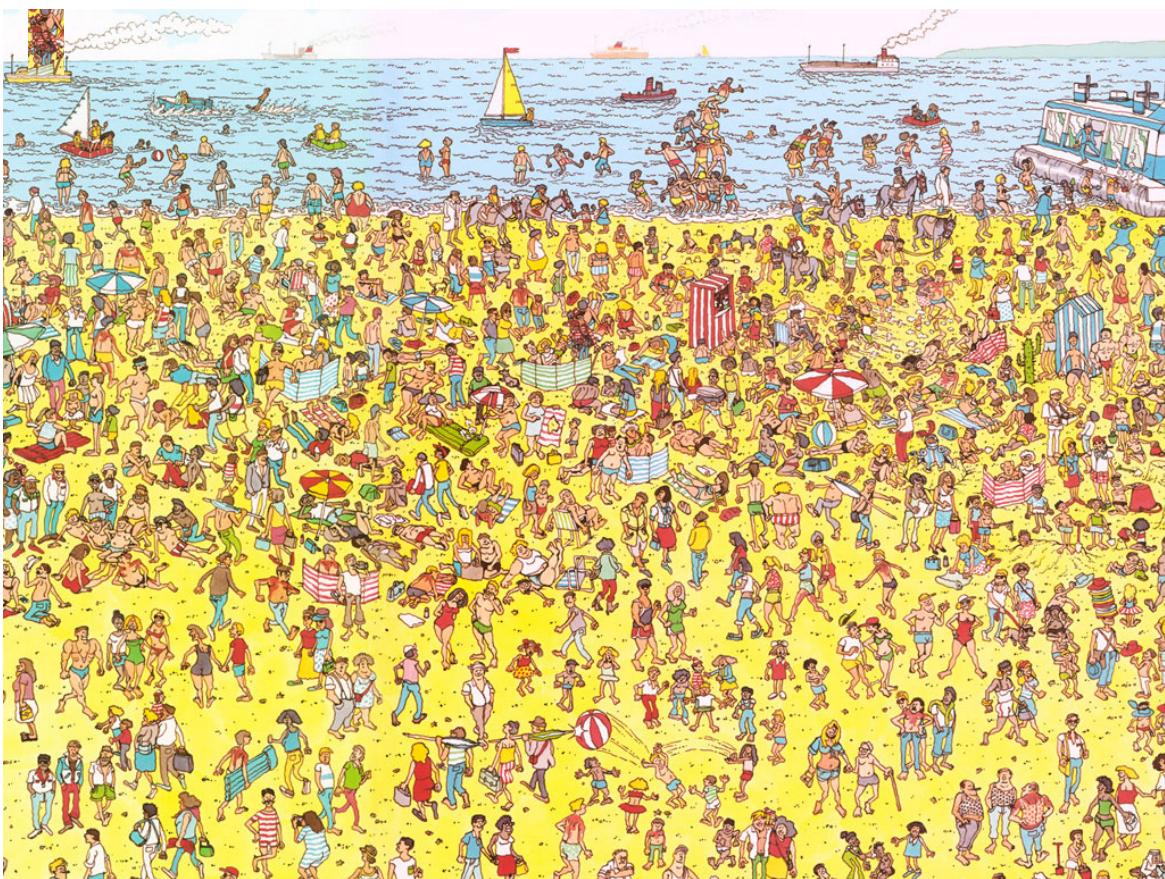
Nam liber tempor cum soluta nobis eleifend option congue nihil imperdiet doming id quod mazim placerat facer possim assum. **Lorem ipsum** dolor sit amet, consectetur adipiscing elit, sed diam nonummy nibh euismod tincidunt ut laoreet dolore magna aliquam erat volutpat. Ut wisi enim ad minim veniam, quis nostrud exerci tation ullamcorper suscipit lobortis nisl ut aliquip ex ea commodo consequat.

Duis autem vel eum iriure dolor in hendrerit in vulputate velit esse molestie consequat, vel illum dolore eu feugiat nulla facilisis.











EIN BEISPIEL: π PARALLEL BERECHNEN

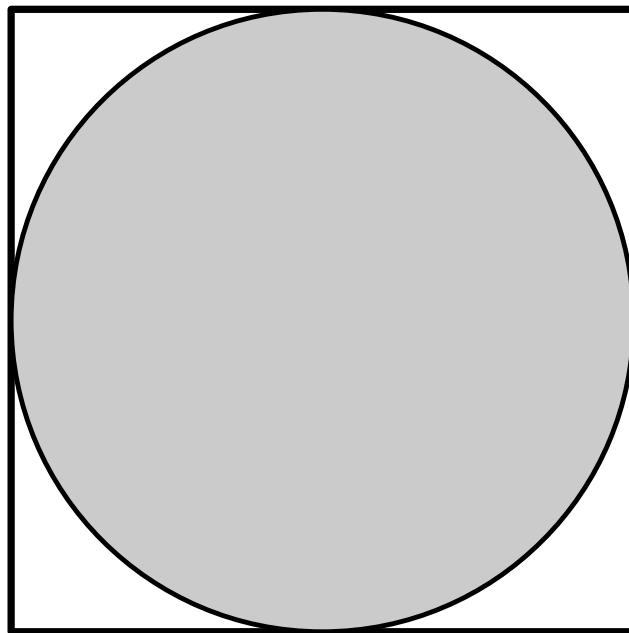


TECHNISCHE
UNIVERSITÄT
DRESDEN

Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 139



$$\text{Fläche} = \pi r^2$$

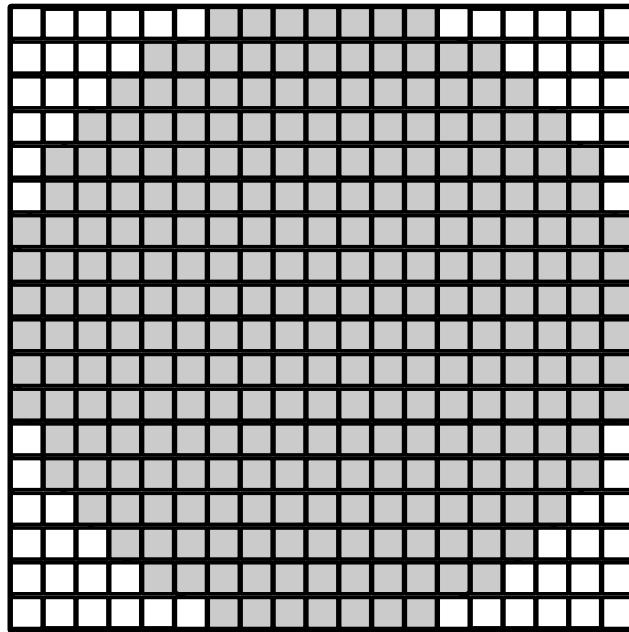


TECHNISCHE
UNIVERSITÄT
DRESDEN

Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 140



Zerlegen der Fläche in Pixel



Muss das Pixel grau oder weiss sein?



```
#define NBLOCKS  = 4096
#define NTHREADS = 128

__global__ void setgreypixels(byte *pixels) {
    int      tidx  = threadIdx.x + blockIdx.x * blockDim.x;           //Thread-ID in
    x
    int      tidy   = threadIdx.y + blockIdx.y * blockDim.y;           //Thread-ID in
    y
    float   x;      = float(tidx) / ( blockDim.x * gridDim.x ) - 0.5; //x-Pos Pixel
    float   y;      = float(tidy) / ( blockDim.y * gridDim.y ) - 0.5; //y-Pos Pixel
    int      pidx  = tidx * blockDim.x * gridDim.x + tidy;            //Pixel-ID

    pixels[pidx] = byte( ( x*x + y*y ) <= ( 0.25 ) ); //Schlecht: C-Konvention!
}

...
dim3 blocksquare(NTHREADS, NTHREADS);
dim3 gridsquare(NBLOCKS, NBLOCKS);

...
setgreypixels<<<gridsquare,blocksquare>>>(devPixels);
```

Pixel zählen und Fläche normieren

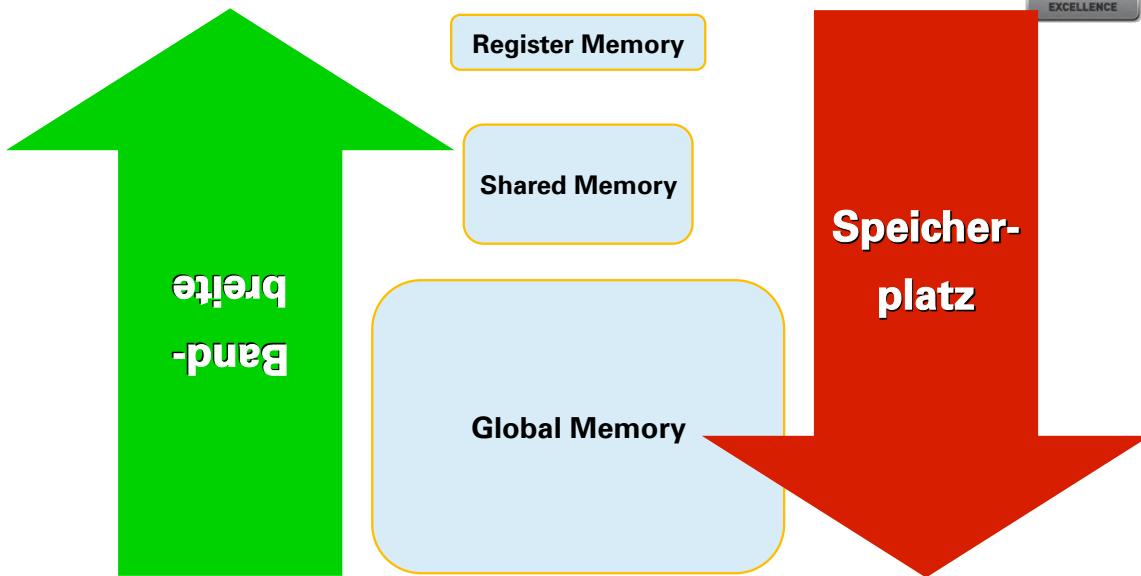


```
int main( int argc, char *argv[] ) {  
  
    ...  
  
    dim3 blocksquare(NTHREADS, NTHREADS);           //quadratischer Block  
    dim3 gridsquare(NBLOCKS, NBLOCKS);                //quadratisches Grid  
  
    int pidx;  
    int npixels = 0;  
  
    float pi;  
  
    ...                                                 //malloc host/device  
  
    setgreypixels<<<gridsquare,blocksquare>>>(devPixels);  
    ...                                                 //cudaMemcpy  
  
    for ( pidx=0; pidx++; pidx<=NTHREADS*NBLOCKS ) npixels += hostPixels[pidx];  
    pi = float( NTHREADS * NBLOCKS ) / float( npixels );  
  
    ...                                                 //printf, free, etc.  
  
}
```

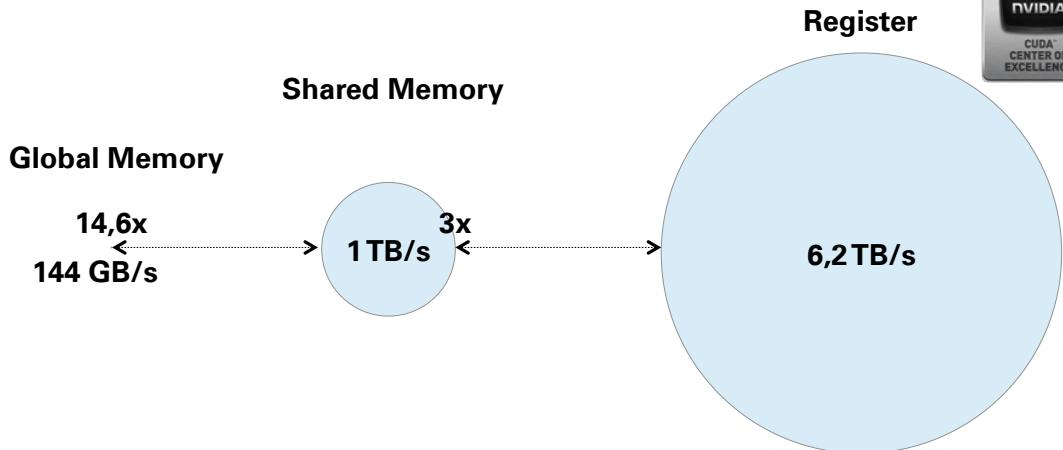


SPEICHERZUGRIFFE

Speicherhierarchie



Speicherbandbreite

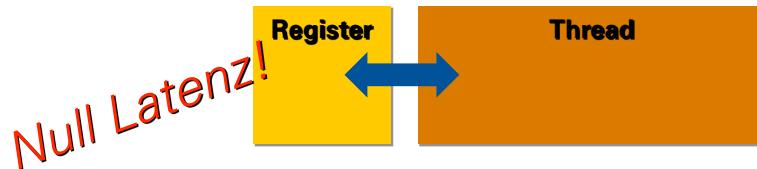


Tesla C2050 (FERMI)

Shared memory: $4 \text{ Byte} * 32 \text{ Banks} * 14 \text{ SM} * 0,5 * 1,15\text{GHz} = 1,03 \text{Tbyte/s}$

Registers (a*b+c): $3 * 4 \text{ Byte} * 14 \text{ SM} * 32 \text{ Threads} * 1.15 \text{ GHz} = 6,18 \text{Tbyte/s}$

Register-Speicher



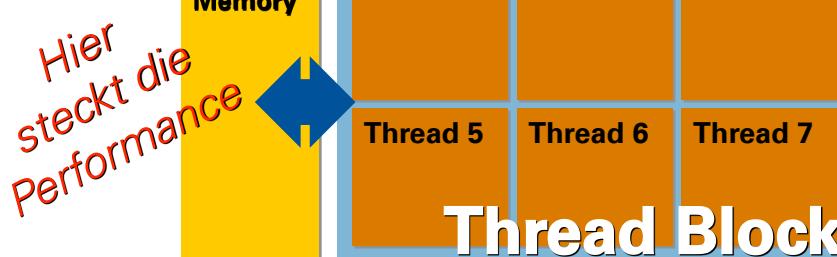
Anzahl von 32 bit Registern pro Multiprozessor

Compute Capability 1.0-1.1: 8k
Compute Capability 1.2-1.3: 16k
Compute Capability 2.x: 32k
Compute Capability 3.0, 3.5: 64k

Was?
Wie?
Wer?
Lebensdauer?
Zugriff?

Registerspeicher ist der lokale Speicher für eine ALU, on-chip
Einfach im Kernel eine Variable deklarieren, z.B. int counter;
Genau einem Thread zugewiesen, nur dieser kann darauf zugreifen
Während der Laufzeit des Threads
Lesen und Schreiben

Shared Memory



Maximaler Shared Memory pro Multiprozessor

Compute Capability 1.0-1.3: 16k
Compute Capability 2.x, 3.0, 3.5: 48k

Was?
Wie?
Wer?
Lebensdauer?
Zugriff?

Speicher zum **schnellen** Datenaustausch für Threads, on-chip
`__shared__ float commondata[threadsPerBlock];`
Alle Threads in einem Block
Während der Laufzeit des Thread Blocks
Lesen und Schreiben



```
__global__ void foo() {
    __shared__ int a[256];
    ...
}
```

Erzeugt ein Integer
Array „a“ mit 256
Elementen im Shared
Memory für jeden
Threadblock

Achtung!
Größe des Arrays im Kernel
festgelegt!

Shared Memory: Dynamische Größe zur Laufzeit



```
__global__ void foo() {
    extern __shared__ int a[];
    ...
}

int main(...){
    ...
    foo<<< NBLOCKS, NTHREADS,
    NTHREADS*sizeof(int)>>>();
    ...
}
```

Erzeugt ein
Integer Array „a“
mit NTHREADS
Elementen im
Shared Memory
für jeden
Threadblock

Shared Memory: Füllen



```
__global__ void foo(int *b_gpu){  
    __shared__ int a[256];  
    idx=... //get global index  
    ...  
    // every thread copies one b to a  
    a[threadIdx.x]=b[idx];  
  
    // wait for all threads in a block  
    __syncthreads;  
    ...  
}
```

Shared Memory: Typische Verwendung

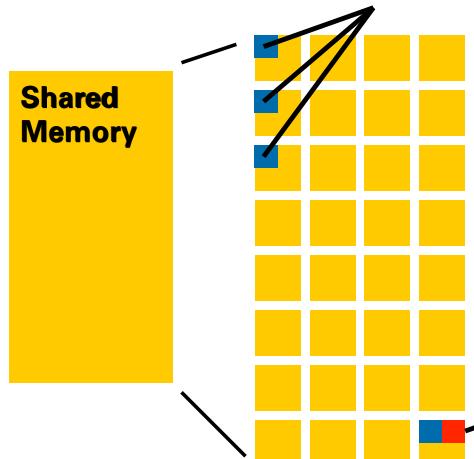


```
__global__ void foo(int *b_gpu){  
    __shared__ int a[256];  
    idx=... //get global index  
    ...  
    // copy reused date to shared memory  
    // sync  
    // compute  
    // sync  
    // write back to global memory  
    // evtl. sync  
}
```

Struktur des Shared Memory: Banks



Hintereinanderliegende 32 bit Worte werden
in hintereinanderliegende banks abgelegt (*interleaving*)

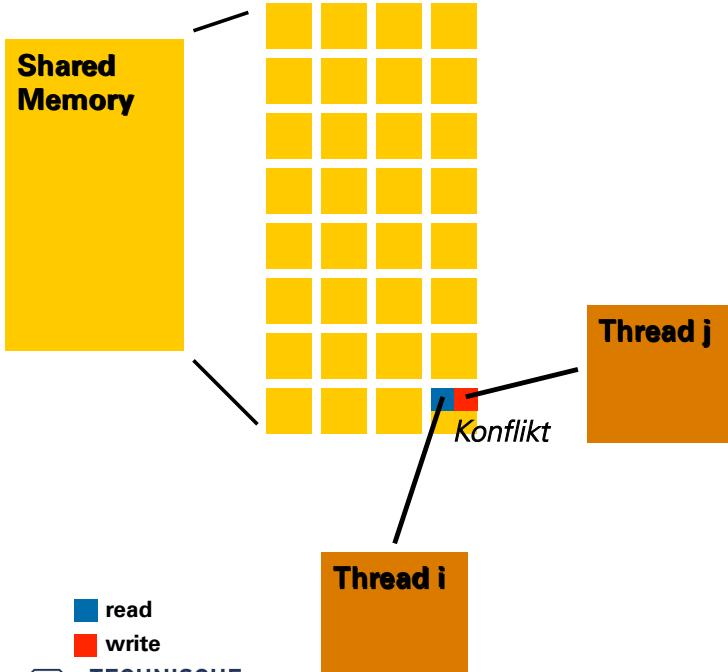


Anzahl der Shared Memory Banks

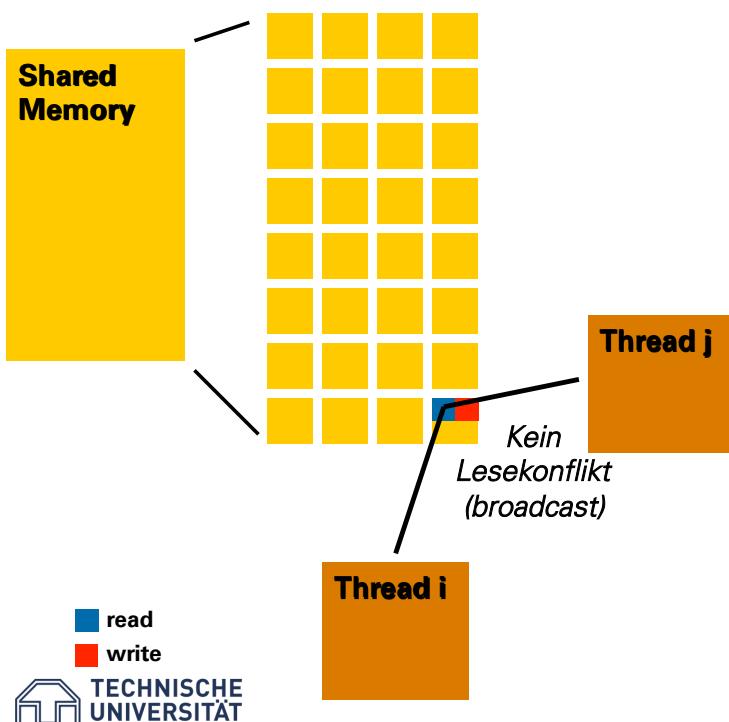
Compute Capability 1.0-1.3: 16

Compute Capability 2.x, 3.0, 3.5: 32

Shared Memory: Bank-Konflikte



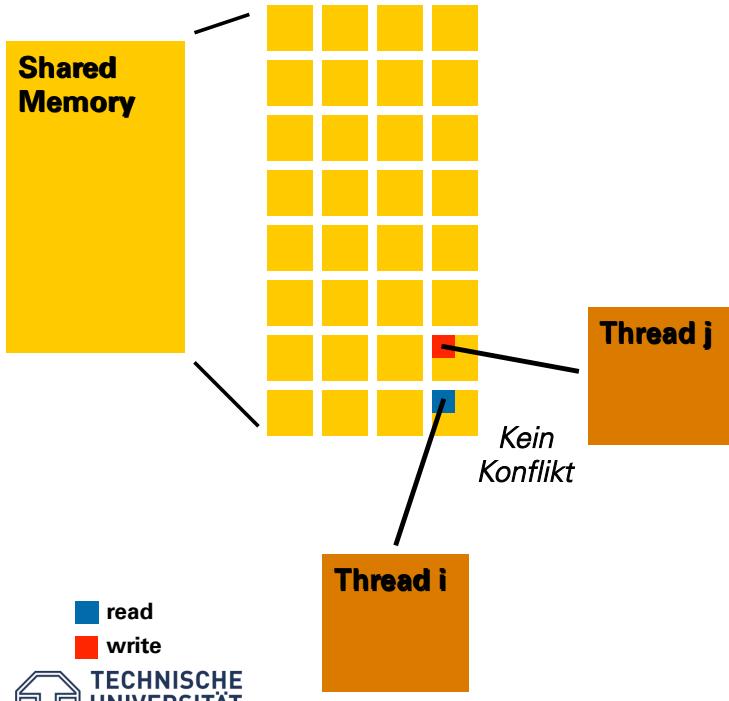
Shared Memory: Bank-Konflikte



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 155



Shared Memory: Bank-Konflikte



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 156



Strided Zugriff



32 Bit

```
__shared__ char shared[32];  
char data = shared[BaseIndex + tid];
```

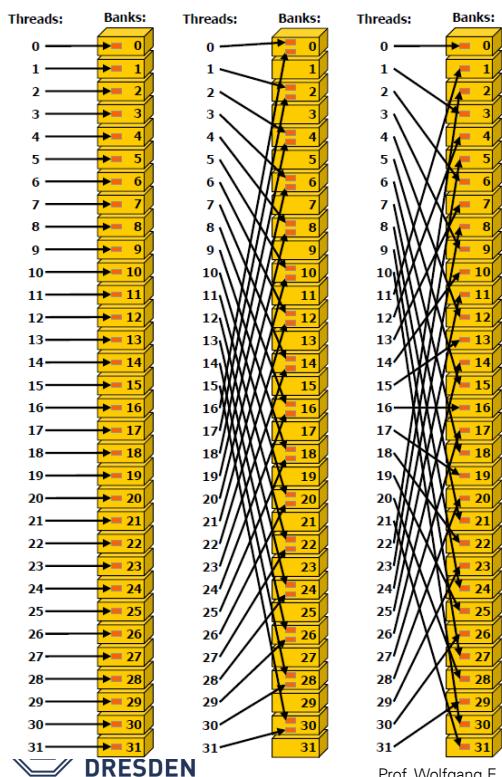
tid und tid+n greifen auf dieselbe Bank zu, falls
 $s*n = \text{Anzahl Banks}$

>32 Bit

```
struct type {  
    float x, y, z;  
};  
  
__shared__ struct type shared[32];  
struct type data = shared[BaseIndex + tid];
```

3 separate 32-Bit-Lesezugriffe ohne Bank-Konflikte,
da auf x, y, z mit einem Stride von 32 Bit zugegriffen wird

Strided Zugriff und Bank-Konflikte (SM 2.x)



Links:

Lineare Adressierung mit 1x32 Bit Stride
Kein Bank-Konflikt

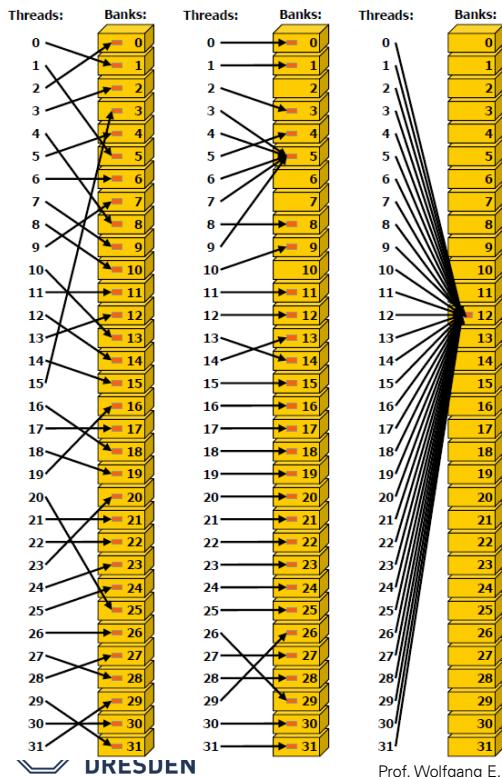
Mitte:

Lineare Adressierung mit 2x32 Bit Stride
2-Wege Bank-Konflikt

Rechts:

Lineare Adressierung mit 3x32 Bit Stride
Kein Bank-Konflikt

Strided Zugriff und Bank-Konflikte (SM 2.x)



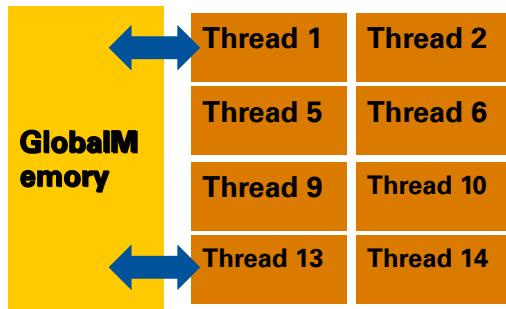
Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 159



Global Memory: Viel Platz, aber laaaaangsaaaaam



Das
kann man
maximal
auf eine
GPU
packen



Derzeit ca. 4 bis 6 GB

Wenn man nur
Global Memory nutzt,
sollte man sich
lieber CPUs kaufen

Was?
Wie?
Wer?
Lebensdauer?
Zugriff?

Speicher zum **langsamen** Datenaustausch zwischen Blocks
`_global_ float commandData[maxMem];`

Alle Threads auf dem Device, Host

Laufzeit des Programms

Lesen und Schreiben

Nicht ausgerichtet gespeicherte Daten (global mem)



Speicher (4 x 128 Byte)



Angeforderter Speicher (1 x 128 Byte not aligned)



Übertragener Speicher (2 x 128 Byte)



50% verschenkte Speicherbandbreite

Ausgerichtet gespeicherte Daten (global memory)



Speicher (4 x 128 Byte)



Angeforderter Speicher (1 x 128 Byte aligned)



Übertragener Speicher (1 x 128 Byte)



2D Speicherzugriff

```
T* elem = (T*) ((char*)Base + Row * pitch) + Column
```

Daten ausgerichtet speichern

Speicher wird ausgerichtet erstellt (auf der Host Seite)



1D Array

- `cudaMalloc()` Startadresse ist ausgerichtet

2D Array

- `cudaMallocPitch()` Startadresse jeder Zeile ist ausgerichtet

Weitere Mallocs für Device Speicher

- `cudaMalloc3D()`, `cudaMallocArray()`,
`cudaMalloc3DArray()`



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 163



Beispiel: Ausrichtung von Strukturen und Klassen

```
struct Foo {  
    float a;  
    float b;  
    float c; }  
// sizeof(Foo) == 12
```

Speicher (16 x 4 Byte)



```
struct __align__(16) Bar {  
    float a;  
    float b;  
    float c; }  
// sizeof(Bar) == 16
```

Speicher (16 x 4 Byte)



Prof. Wolfgang E. Nagel / Guido Juckeland – Slide 164

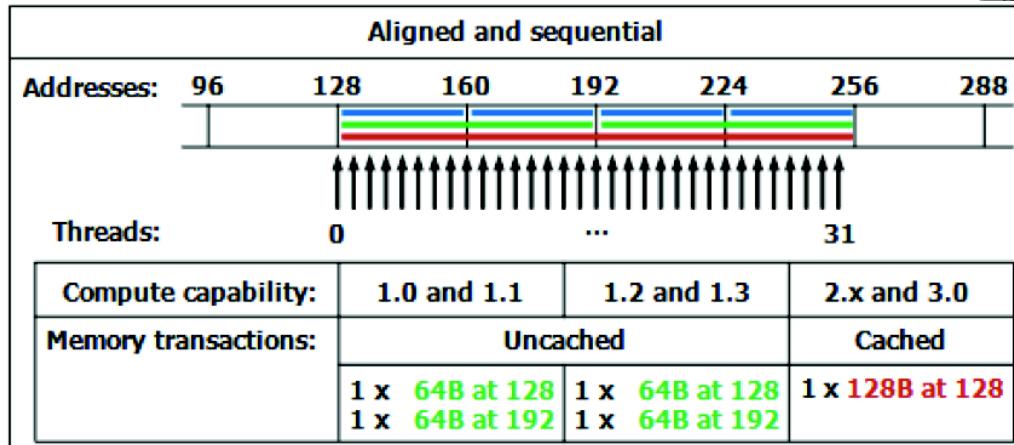


Caching für Global Memory (ab CC 2.x)



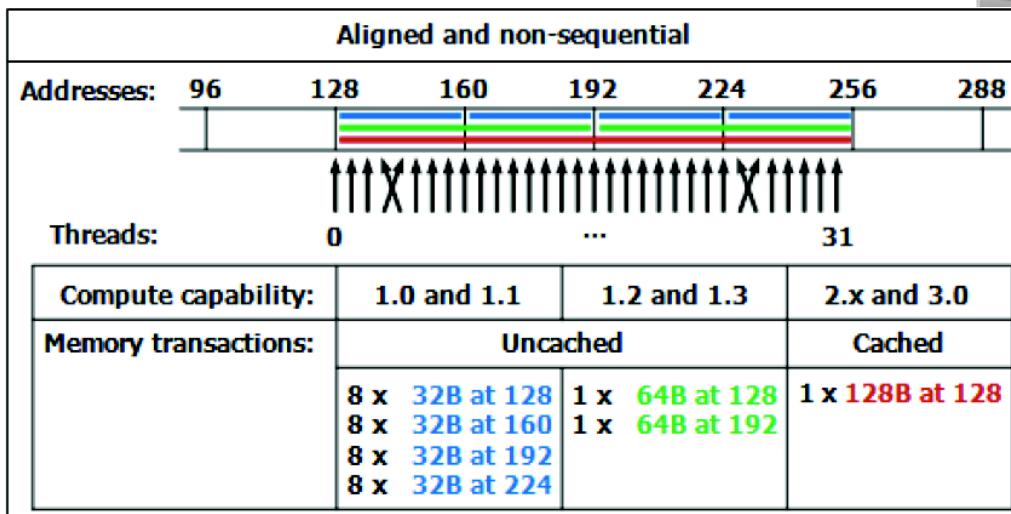
- Wie?** Compiler-Flag `-dlcm`
`-xptxas -dlcm=ca` (L1 und L2 Cache genutzt)
`-xptxas -dlcm=cg` (nur L2 Cache genutzt)
- Was passiert?** Es werden immer 128 Bytes gecached (CC 2.x)
Diese entsprechen ausgerichteten 128 Byte-Segmenten
Mit L1+L2 Cache werden einmal 128 Byte transferiert
Mit L2 Cache alleine werden 4x32 Byte transferiert
- Und dann?** Wird mehr als 128 Byte pro Thread benötigt, so wird eine Speicheranfrage pro Warp in separate 128-Byte-Anfragen gesplittet
- Parallelität?** 128-Byte-Anfragen werden unabhängig voneinander bearbeitet
Ist das Datum im Cache, geschieht dies mit der Cache-Geschwindigkeit, sonst mit der Geschwindigkeit des Global Memory

Ausgerichteter, sequentieller Zugriff auf Global Mem.



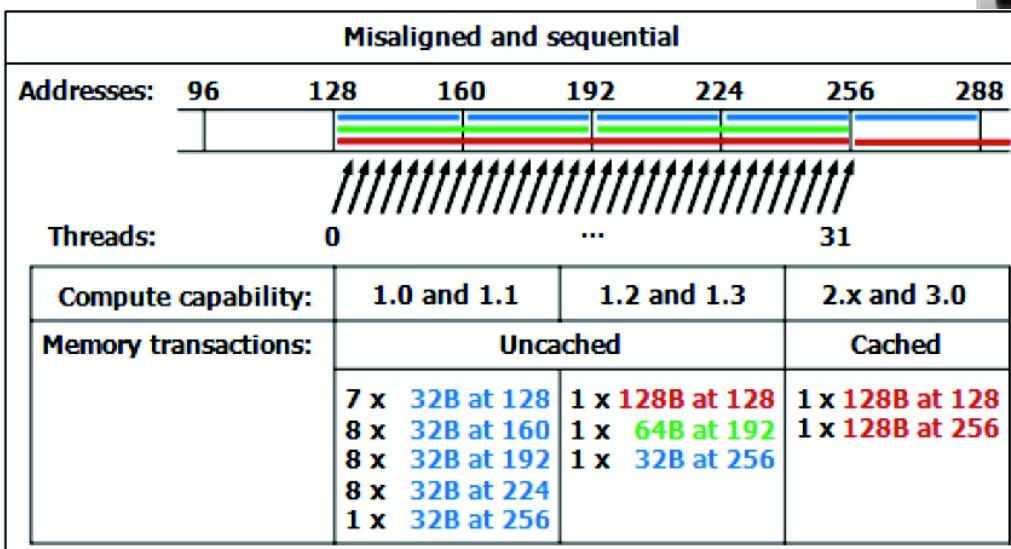
Warp-Zugriff, ein 4 Byte Wort pro Thread

Ausgerichteter, nicht sequentieller Zugriff



Warp-Zugriff, ein 4 Byte Wort pro Thread

Nicht ausgerichteter, sequentieller Zugriff



Warp-Zugriff, ein 4 Byte Wort pro Thread



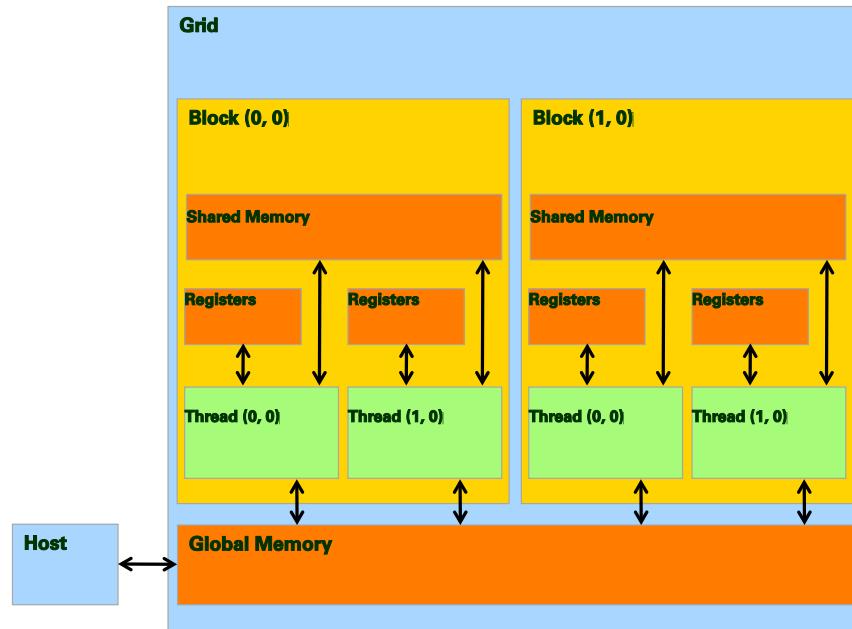
Vorteile

- Ausnutzung theoretischer Speicherbandbreite
- Bessere Ausnutzung L1/L2 Cache
- Weniger Kollisionen beim Speicherzugriff
- Vermeidet Nutzung von Local Memory durch nvcc

Nachteile

- Umständliche Speicherzugriffe
- Eventuell langsame 2D/3D cudaMemcpy host2device/ device2host

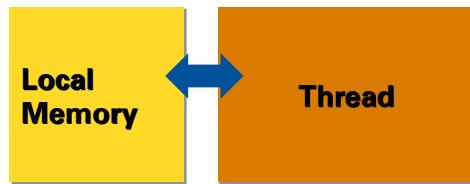
Übersicht der Hardware-Speicher inklusive Zugriff



Local Memory: Getarnter Global Memory



Kann
Concurrency
verhindern!
Do not use!



Alles, was nicht
in Register passt
oder keine definierte
Größe hat

Größe des Local Memory pro Thread

Compute Capability 1.0-1.3: 16kB
Compute Capability 2.x: 512kB

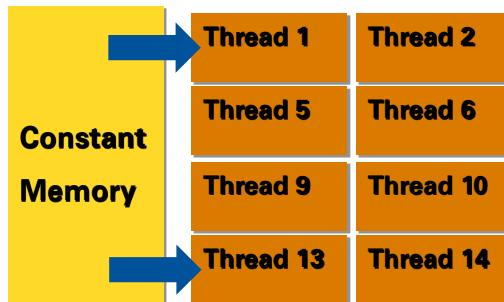
Was?
Wie?
Wer?
Lebensdauer?
Zugriff?

Globaler Memory, Pro 32 bit Wort genau ein Thread, coalesced
`__local__ int counter;`
Genau einem Thread zugewiesen, nur dieser kann darauf zugreifen
Während der Laufzeit des Threads
Lesen und Schreiben

Constant Memory: Read-Only Device Memory



Kann
Bandbreite
sparen



Wenn Du es
nicht ändern mußt,
geht's schneller.

Größe des Constant Memory

Compute Capability 1.0-2.x: 64kB

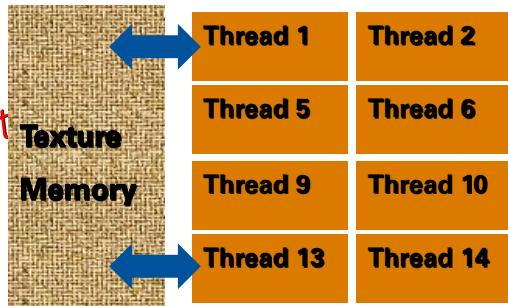
Was?
Wie?
Wer?
Lebensdauer?
Zugriff?

Device Memory, gecashed (8kB), Broadcast (spart Lesezugriffe)
`__constant__ float ReadOnlyData [N];`
Alle Threads auf dem Device, Host
Laufzeit des Programms
Nur Lesen

Texture Memory: Grafischer Speicher für Nostalgiker



Cached,
2D-optimiert



Ist fast wie
OpenGL oder
Direct3D

Derzeit ca. 4 bis 6 GB

Was?

Im Grunde Global Memory, zugriffsoptimiert, OpenGL/Direct3D

Wie?

```
texture<DataType, TextureType, ReadMode> tex;  
cudaBindTextureToArray(tex, cudaArray);
```

Wer?

Alle Threads auf dem Device, Host

Lebensdauer?

Laufzeit des Programms

Zugriff?

Nur Lesen