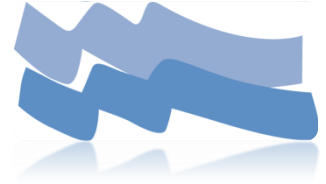


# **Potenzial der NVIDIA Graphik Prozessoren (GPU) und der INTEL XEON Phi Coprozessor-Technik**

**Dr.- Ing. Hartmut Sacher und Dr. - Ing. Alpaslan Yörük**

# Gliederung

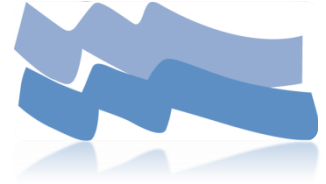


- Aktuelle und künftige Parallelisierung in Hydro\_AS-2D
- Benötigte Hardware
- Vergleich der Simulationsdauer für verschiedene Berechnungsbeispiele  
Hydro\_AS-2D 3.16 vs. Hydro\_AS-2D GPU
- Hinweise und Fazit

# Parallelisierung – Grundbegriffe und Hinweise



- Aktuelle Parallelisierung in Hydro\_AS-2D 3.15
  - Erfolgt durch Compiler Optionen, die Intel Multicore-Architekturen unterstützen.
  - D.h., es werden mehrere Prozessoren für die Berechnung parallel genutzt.
- Parallelisierung in Hydro\_AS-2D GPU
  - GPU-Parallelisierung – Unterstützung von Graphikprozessoren (NVIDIA)
  - Coprozessoren – Unterstützung von numerischen Coprozessoren (INTEL)
  - Eine Anpassung der Software (Berechnungscode) ist erforderlich.
  - Die Umsetzung wird gemacht, da die Parallelisierung gut machbar ist und deutliche Verkürzungen der Simulationsdauer möglich ist.
  - Die Aufgabe ist aber komplex und wird kontinuierlich weiter verfolgt.
- Anwenderhinweise
  - Es ist eine spezielle Hardware erforderlich.
  - Für den Coprozessor-Einsatz ist z.Z. ein Linux Betriebssystem erforderlich.
  - Es wird weiter eine „normale“ Hydro\_AS-2D Version geben.



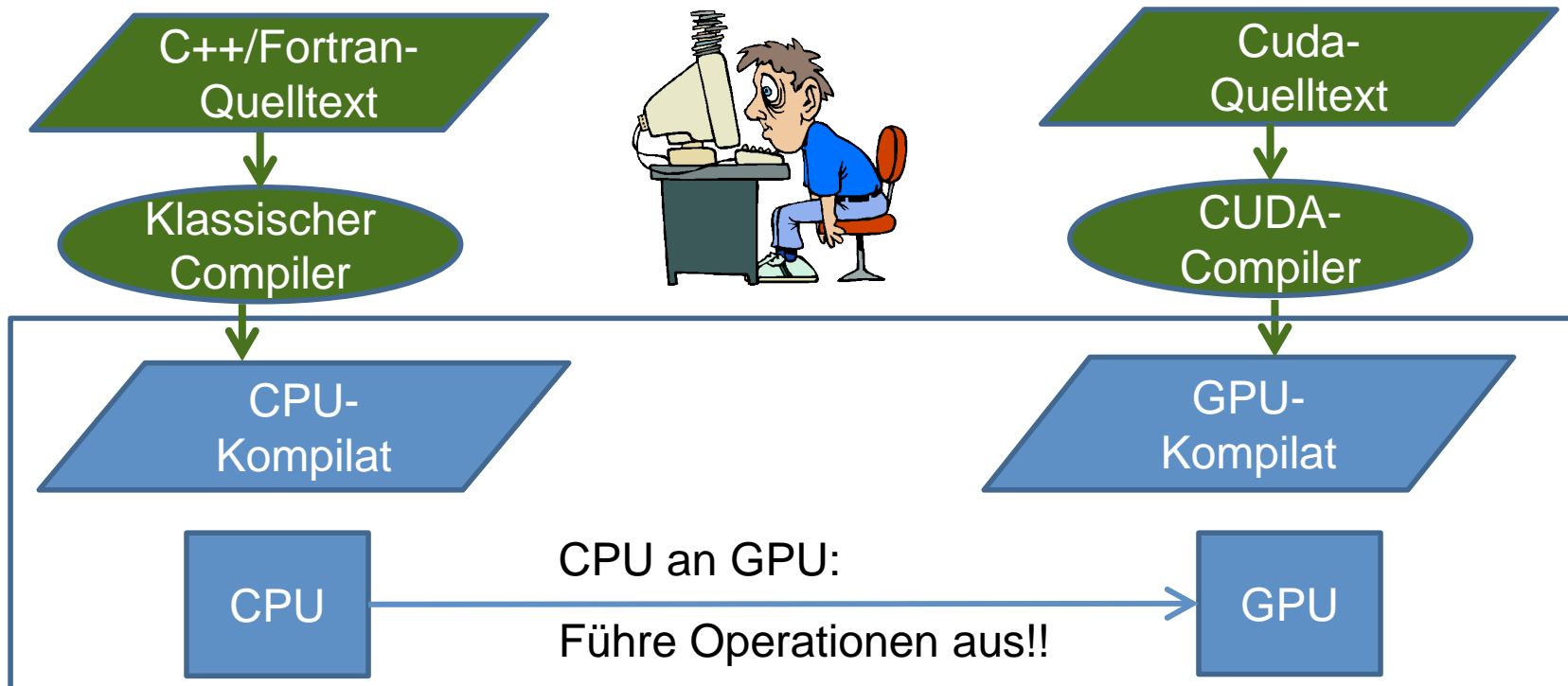
# Parallelitätsmodelle

- MIMD: Multiple Instruction Multiple Data
  - Unabhängig agierende Prozessoren
  - Beste Auslastung bei unabhängigen Aufgaben
  - Bei Abhängigkeiten in Instruktionen: Synchronisation, Mutex (mutual exclusion),...
  - Realisiert in heutigen Rechnern mit mehreren Cores
  - Und seit ca. 2013 ... auch in INTEL XEON Phi Coprozessoren
  
- SIMD: Single Instruction Multiple Data
  - Vektoroperationen:  $A[1:100000] = B[1:100000] + C[1:100000]$
  - Beste Auslastung bei identischen Aufgaben
  - Parallelisiert „klassisches“ `for i= 1 .. n`
  - Berechnung durch Einsatz von Graphikprozessoren (GPU, NVIDIA)
  - Problem: `if-then-else` ist hochgeradig *nicht* SIMD;  
d.h., der GPU-Ansatz ist in Wahrheit komplex...

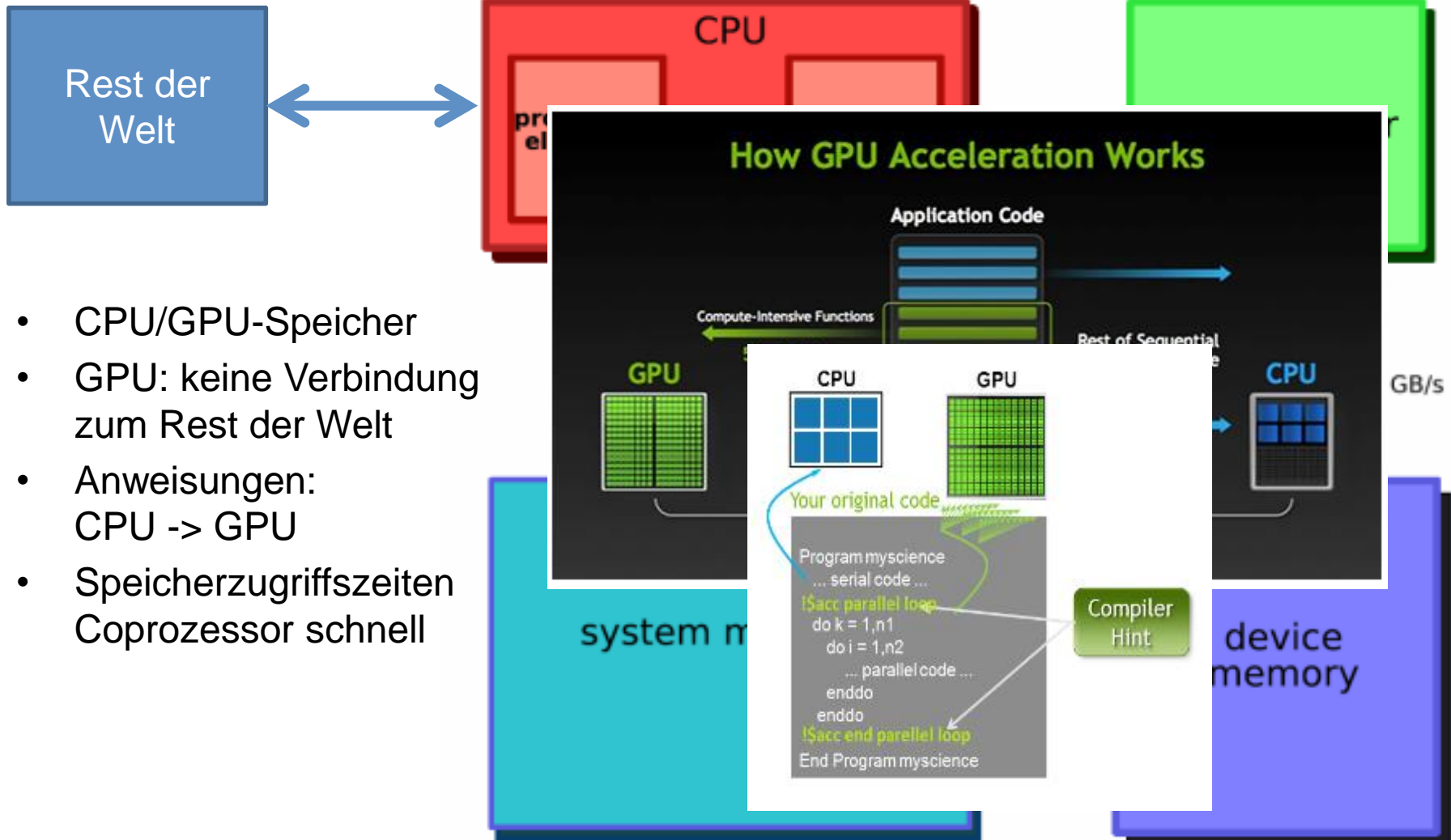
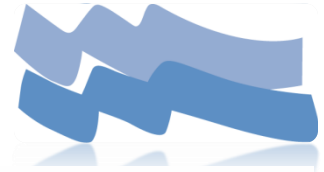
# Wie programmiert man GPU?



- Spezielle Programmiersprachen
  - NVIDIA: CUDA (CUDA: Compute Unified Device Architecture) oder herstellerunabhängig: OpenCL
- GPU-Programmteile werden von CPU-Programm verwaltet und an die GPU „zur Berechnung“ übergeben



# Wie funktionieren GPU im Rechner?



- CPU/GPU-Speicher
- GPU: keine Verbindung zum Rest der Welt
- Anweisungen: CPU -> GPU
- Speicherzugriffszeiten Coprozessor schnell

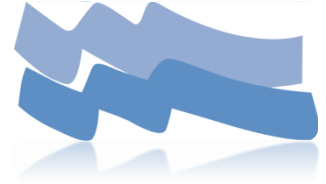
# Simulations-Systeme unserer Tests



- Simulationsrechner Intel CPU (2x QuadCore)
- NVIDIA GPU Quadro K2000 (ca. 400 €)
- NVIDIA GPU Tesla K20m (ab ca. 3000 €)
- NVIDIA GPU Tesla K40 (ab ca. 5000 €)
- Unsere Testumgebung (4xK20): ca. 13.000 €
  
- Es sind auch mehrere GPU-Einheiten in einem Rechner möglich. Aber nur ein Rechenlauf pro GPU.



<b>CPU-Simulation</b>	<b>GPU Quadro K2000</b>	<b>GPU Tesla K20</b>	<b>GPU Tesla K40</b>
2x Intel® XEON® CPU E5620 @ 2.40GHz	1x Intel® XEON® CPU W3520 @ 2.67GHz	2x Xeon E5-2609 v2 @ 2.50GHz	2x Xeon E5-2660 v2 @ 2.50 GHz
8 GB RAM	2 GB RAM	5 GB RAM	12GB RAM
8 Cores	384 „Cores“	2496 „Cores“	2880 „Cores“



# Berechnungsbeispiele - Modelle

- **Modell 1:** Dambruch-Modell
  - A: 148.000 Knoten und Elemente
  - B: 590.000 Knoten und Elemente
- **Modell 2:** Kleineres Gewässer
  - 670.000 Knoten und 1,1 Mio. Elemente
- **Modell 3:** Gewässersystem (Rheinauen)
  - 1,35 Mio. Knoten und 2,63 Mio. Elemente (instationär)
- **Modell 4:** Großes Gewässer (Donauabschnitt)
  - 570.000 Knoten und 1,1 Mio. Elemente

	<b>Nodes</b>	<b>Modellzeiten</b>	<b>Berechnungszeit CPU</b>
<b>Modell 1A</b>	148185	17 Min	3 h
<b>Modell 1B</b>	589265	17 Min	15 h
<b>Modell 2</b>	668096	28 h	23 d
<b>Modell 3</b>	1349706	111 h	24 d
<b>Modell 4</b>	574751	17 h	26 h



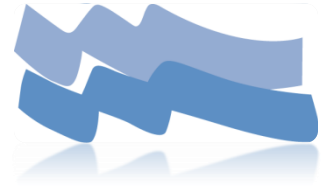
# Berechnungsbeispiele – Beschleunigung GPU



- Deutliche Beschleunigung der Rechenzeiten bei GPU Einsatz im Vergleich zu INTEL CPU
- Beschleunigung nimmt mit Anzahl GPU-Cores zu (Höhere Anzahl Graphikprozessoren bzw. Systemleistung)
- Beschleunigungsfaktor ist für alle Berechnungsmodelle unterschiedlich!

				Beschleunigungsfaktor		
	Knoten	Modellzeiten	Berechnungszeit CPU	Quadro K2000	Tesla K20	Tesla K40
<b>Modell 1A</b>	148185	17 Min	3 h	20	35	41
<b>Modell 1B</b>	589265	17 Min	15 h	17	36	42
<b>Modell 2</b>	668096	28 h	23 d	6	13	17
<b>Modell 3</b>	1349706	111 h	24 d	5	11	12
<b>Modell 4</b>	574751	17 h	4 d	5	20	25

# Berechnungsbeispiele – Beschleunigung GPU



- Beschleunigung bei GPU Einsatz nimmt mit **Anteil benetzter Elemente** zu

Test-Modelle	Anteil benetzter Elemente [%]	Beschleunigungsfaktor		
		K2000	K20	K40
Modell 1A	95	20	35	41
Modell 1B	95	17	36	42
Modell 2	54	6	13	17
Modell 3	25	5	11	12
Modell 4	30	7	19	21

- Zusätzliche Rechenläufe Modell 4 im Vergleich zur CPU-Berechnung

Weitere Tests Modell 4	Anteil benetzter Elemente [%]	Berechnungszeit CPU	Beschleunigungsfaktor		
			K2000	K20	K40
Modell 4a HQ100	30	4 d	5	20	25
Modell 4b MQ	5	1 d	5	19	21
Modell 4c HQextrem	52	6 d	7	26	30

# INTEL Xeon Phi Coprozessor



- Der XEON Prozessor ist der „Host“ für den Coprozessor, der mit der MIC (Many Integrated Core) Technologie arbeitet.
- Die Kommunikation läuft über den PCIe Bus, der Coprozessor läuft unter LINUX und kommuniziert mit dem PCIe Bus via TCP/IP.
- Ein Simulationsprogramm läuft also sowohl auf dem Host als auch als individueller Job auf dem Coprozessor. D.h., ein Coprozessor kann gleichzeitig mehrere Jobs verarbeiten.

## Intel® Xeon Phi™ Coprozessoren



Die Intel Xeon Phi Coprozessoren bieten bis zu 61 Kerne, 244 Threads und 1,2 TeraFLOPS Leistung. Sie sind in verschiedenen Konfigurationen für unterschiedliche Hardware-, Software-, Workload-, Leistungs- und Effizianzorderungen verfügbar.

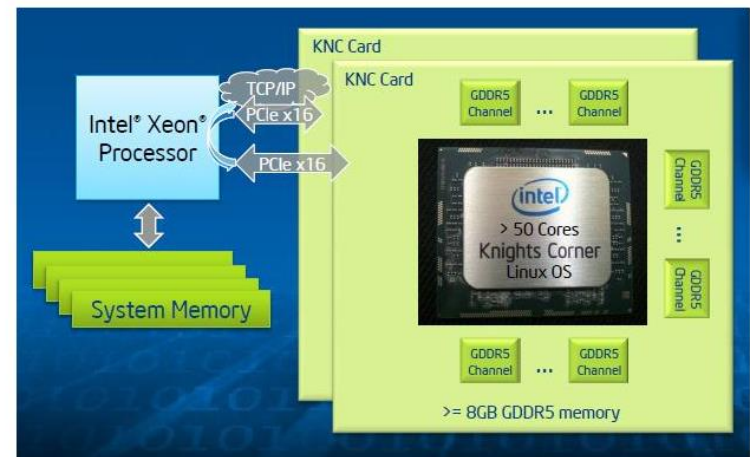


Figure 1. The first generation Intel® Xeon Phi™ product codenamed “Knights Corner”

# INTEL Xeon Phi Coprozessor

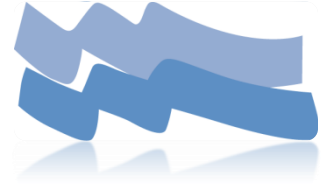


- Unsere Testumgebung (ca. 4.500 €):

<b>CPU-Simulation</b>	<b>INTEL XEON Phi Coprozessor</b>
2x Intel® XEON® CPU E5-2609 @ 2.50GHz	1xINTEL Xeon Phi 3120P, 1.100 GHz
32 GB RAM	6 GB RAM
8 Cores	57 „Cores“

- 61 Kerne nicht gekoppelt mit der CPU, jeder Kern ist sehr leistungsfähig (x-faches von dem, was ein GPU-CUDA-Kern leistet).
- Die Software muss für die Coprozessor-Unterstützung angepasst werden.
- Daran arbeiten wir zur Zeit. D.h., wir können noch keine Aussage über die Simulationsdauern machen.

# Bisheriges Fazit



- Der Einsatz einer GPU bringt deutliche Reduzierung der Simulationszeiten.
- Mehrere GPU in einem Rechner sind möglich. Mehrere Rechenläufe auf einer GPU sind nicht sinnvoll. Effizient ist z.Z. nur ein Rechenlauf je GPU.
- INTEL Coprozessor Technik werden wir weiter verfolgen; ggfs. sind je nach Anwendungsfall GPU oder Coprozessor schneller.
- Leistungsfähigere Hardware ist zu verwenden.
- 10 Mio. Knoten Netze? Dann muss ggfs. ein Gesamtmodell in einzelne Gebiete aufgeteilt werden (Performance SMS beachten).
- Ggfs. Umring testen für optimale Gebietsgröße.



# **Potenzial der NVIDIA Graphik Prozessoren (GPU) und der INTEL XEON Phi Coprozessor-Technik**

**Dr.- Ing. Hartmut Sacher und Dr. - Ing. Alpaslan Yörük**