



TECHNISCHE
UNIVERSITÄT
DRESDEN

DIGITAL SIGNAL PROCESSING USING CUDA

Digital Signal Processing using CUDA

Nico Wehmeyer, Richard Pfeifer, Fabian Jung

Dresden, 3.2.2014



DRESDEN
concept
Erleben des
Wissenschafts
und Kultur

Inhalt

Aufgabenstellung

Überblick

Host Code

Levenberg Marquardt

Verwaltung Devices

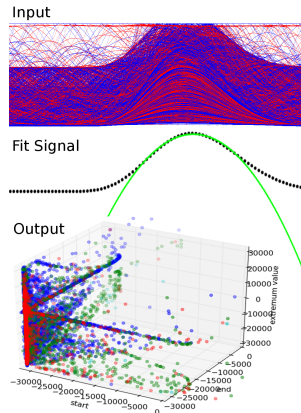
Benchmark

Skalierbarkeit

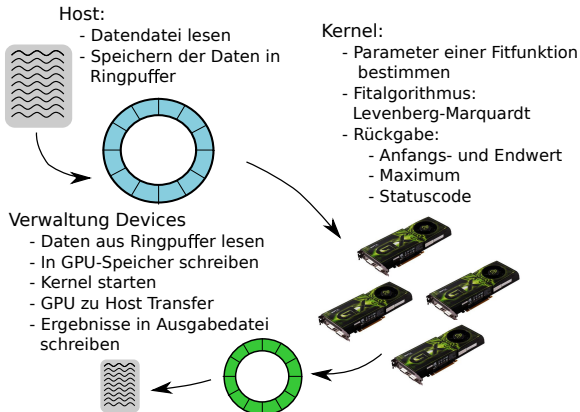
Danksagung

01 Aufgabenstellung

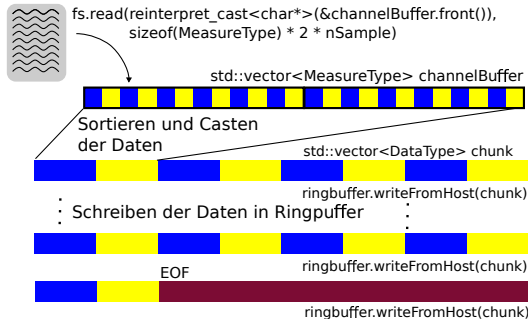
- Messgeräte erzeugen kontinuierlichen Datenstrom
- Daten: Kurven mit Tausenden von Samplepunkten
- Ziel: Datenreduktion durch Fitten von Vorgabefunktion
- Anforderungen:
 - Hoher Datendurchsatz
 - Skalierbar auf bis zu 4 GPUs/Node



Überblick



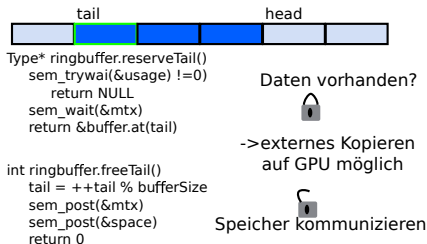
Host: Lesen der Daten



- Lesen von Daten zweier Kanäle
- Kanaltrennung
- Casten der Daten
- Sammeln mehrerer Signale zu Chunks
- Letzter Chunk wird mit Nullen gefüllt.

Host: Ringpuffer

- Messdatenfile → DataReader → Ringpuffer → GPU
- Speicher des Ringpuffer: Vektor dessen Typ per Template frei gewählt werden kann
- Host-GPU-Transfer wird wie folgt gelöst:



Levenberg Marquardt (1)

- Eingabedaten
 - Samples
 - Compute Capability 1.x: ca. 800)
 - Compute Capability 2.0 oder höher: ca. 2500)
 - Interpolationsschritt
 - beliebige Dezimalzahl größer 0)
 - Interpolation durch Texture Memory)

Levenberg Marquardt (2)

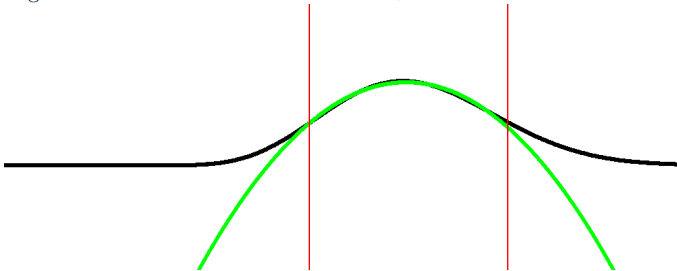
- Verarbeitung
 - eine Ausgleichungsrechnung pro Block
 - Grund: ca. das 5-fache der Sampleanzahl an Shared Memory benötigt (bei 1000 Samples ca. 20 kB)
 - Zugriff auf Samples durch Texture Memory
 - Shared Memory gespart
 - schnelle Interpolation möglich
 - Vorgehensweise
 - Anfangs- und Endwert ermitteln
 - abhängig vom Schwellwert Bereich festlegen
 - für den Bereich Näherungsfunktion ermitteln
 - Qualität durch Residuen und Maximum der Funktion ermitteln

Levenberg Marquardt (3)

- Ausgabedaten
 - 3 Parameter einer quadratischen Funktion: $a * x^2 + b * x + c$
 - Anfangs- und Endwert
 - Maximum
 - durchschnittliche Abweichung
 - Status (Fehler, Erfolg, Abbruch)

Levenberg Marquardt (4)

- Ergebnis: $\text{fitfunktion} = -0.550151 * x^2 + 654.15509 * x - 203167.921875$



Verwaltung Devices

- Jedes Devices wird von einem Thread verwaltet
- Asynchrone Aufrufe
 - `cudaMemcpyAsync`
 - Kernel
 - Ermöglicht Pipeline
- Eigener Thread für Ausgabedatenstrom

Benchmark

GPUs	Datei	Zeit	Rate	Normiert
1	247MB	72.450s	3.41 MB/s	3.41 MB/s
2	247MB	38.141s	6.48 MB/s	3.24 MB/s
3	247MB	29.060s	8.50 MB/s	2.83 MB/s
4	247MB	20.828s	11.85 MB/s	2.96 MB/s
1	382MB	107.902s	3.54 MB/s	3.54 MB/s
2	382MB	57.537s	6.64 MB/s	3.32 MB/s
3	382MB	40.951s	9.33 MB/s	3.11 MB/s
4	382MB	29.885s	12.78 MB/s	3.20 MB/s

Skalierbarkeit

- Horizontal Skalierbar
- PCIe Bus noch nicht vollständig ausgelastet
- Reserven in der Implementierung des Ringbuffers

Danksagung

Vielen Dank an
Dr. Michael Bussmann, Axel Hübel und Rene Widera
für Diskussion und Optimierungsvorschläge.