

# Mehrstufige zeit- und bewegungsabhängige Rauschreduktion in Echtzeit mittels CUDA

Elmar Bührlé<sup>1</sup>, Benjamin Keck<sup>1,2</sup>, Stefan Böhm<sup>3</sup> and Joachim Hornegger<sup>1</sup>

<sup>1</sup>Lehrstuhl für Mustererkennung, Martensstraße 3, 91058 Erlangen

<sup>2</sup>Siemens Healthcare, CV, Medical Electronics & Imaging Solutions,  
Postfach 3260, 91050 Erlangen

<sup>3</sup>Siemens Healthcare, Imaging & IT Angiography, Siemensstr. 1, 91301 Forchheim

## Zusammenfassung.

## 1 Einleitung

Die Bildgebung interventionell genutzter Angiographieanlagen mittels Röntgentechnologie stellt heute sehr hohe Anforderungen an die Bildqualität. Diese wird durch verschiedenste Faktoren beeinflusst und ist für den Echtzeitbetrieb sicherzustellen, d.h. typischerweise sind die Algorithmen zur Bildverbesserung für einen maximalen Datenstrom von 480-600MBits/s (monoplan-System) bei einer maximalen Latenzzeit von ca. 100 ms auszuführen. Für unsere Untersuchung haben wir einen mehrstufigen zeit- und bewegungsabhängigen Rauschreduktionsalgorithmus gewählt, der bereits in Angiographiegeräten verwendet wird. Die Berechnungsgeschwindigkeit von Bildverarbeitungs-algorithmen ist sowohl von der Komplexität der Algorithmen als auch von der verwendeten Hardware abhängig. Um die Echtzeitfähigkeit des Systems beim Verarbeiten der akquirierten Bildfolge zu garantieren, wird daher bei vielen Produkten auf Spezialhardware, wie z.B. DSPs oder FPGAs, zurückgegriffen. Allerdings müssen Hersteller dafür einen erhöhten Aufwand für die Entwicklung der Spezialhardware als auch für die Implementierung der Algorithmen auf der Spezialhardware in Kauf nehmen.

Im Jahre 2007 stellte der Grafikprozessorhersteller NVIDIA gemeinsam mit einer neuen Generation von Grafikprozessoren (GPUs) eine neue erweiterte Programmierschnittstelle namens "Common Unified Device Architecture"(CUDA) vor, welche die allgemeine Verwendbarkeit der neuartigen Architektur beschreibt. Mithilfe einfacher Erweiterungen der Programmiersprache C kann eine CUDA-fähige Grafikkarte zur Berechnung verschiedenster Algorithmen verwendet werden. Die erste Generation entsprechender GPUs (G80) erreichte durch die extreme Parallelität bereits eine Leistung von mehr als 300 Gigaflops.

In unserer Untersuchung gehen wir auf die Eignung der im Vergleich zu Spezialhardware günstigeren, sowohl leicht integrierbaren als auch vergleichsweise einfach programmierbaren Hardware am Beispiel des gewählten Algorithmus ein.

## 2 Stand der Forschung und Beitrag

Die Eignung von Grafikkarten zur Beschleunigung von Berechnungen im Bereich der Medizinischen Bildverarbeitung wird bereits seit über einem Jahrzehnt erforscht [1].

Mack et. al. zeigte bereits 2004 [2] die Verwendbarkeit von Grafikkarten im Bereich der "Echtzeit-Röntgenbildverarbeitung", wobei die verwendeten Algorithmen mittels Direct3D auf die bisherige Rendering Pipeline abgebildet wurden. Algorithmen werden dabei mittels Shader-Programmiersprachen wie GLSL oder HLSL in die spezialisierte Grafikkarte abgebildet. Die einerseits vergleichsweise einfache Umsetzung von Algorithmen auf GPUs mittels CUDA und die andererseits enorme Rechenleistung wurde bereits 2007 im Bereich der 3D-Rekonstruktion gezeigt [3]. In unserem Beitrag möchten wir die Eignung von GPUs zur Echtzeit 2D-Bildverarbeitung mittels CUDA und die dabei zu beachtenden Programmier-Techniken anhand eines Beispiels erläutern.

### 3 Methoden und Algorithmen

Bei dem beschriebenen Filter handelt es sich um ein rekursives zeitliches Filter erster Ordnung, das die Filterwirkung an die Bewegung im Bild adaptiert. Die Bewegungserkennung und Filterung erfolgt in einer Auflösungshierarchie. Im Folgenden wird das Filter als multiskalares Rauschglättungsfilter mit Bewegungsdetektion bezeichnet. Die wesentlichen Schritte der Filterung sind eine Dekomposition in eine Laplace-Auflösungshierarchie [4] und eine kantenerhaltende Rauschglättung auf variablen Stufen der Laplace-Pyramide.

Vor der Rekonstruktion des Ergebnisbildes wird die Differenz des Stufenbildes zu den vorhergehenden gefilterten Stufenbildern bzgl. der korrespondierenden Stufe gebildet. Aus der niedrigsten Differenz wird ein Gewichtungsfaktor errechnet, der den Einfluss des Bildes mit niedrigster Differenz auf das momentan verarbeitete Bild der entsprechenden Stufe steuert.

#### 3.1 Laplace-Auflösungshierarchie

Die Erstellung einer Laplace-Pyramide gilt als ein wesentlicher Algorithmus in der Bildverarbeitung [4] und wird im beschriebenen Algorithmus zur Dekomposition verwendet. Die anschließende Ortsfilterung dient einerseits als Vorverarbeitungsschritt für die Bewegungsdetektion und andererseits als Ergänzung zum Temporalfilter. Bei der Verwendung von GPUs zur Beschleunigung rechenintensiver Anwendungen muss Folgendes beachtet werden: Teure Zugriffe auf den globalen Grafikspeicher, die jeweils 400 bis 600 Taktzyklen (vgl. [5], S.49) benötigen, können durch den Einsatz von Texturen (Texturcache) vermieden werden. Zugriffe auf den Texturcache erreichen dabei eine ähnliche Latenz wie auf Register. NVIDIAs Speicherverwaltungseinheit unterstützt das Lesen von 32-, 64- oder 128-bit Werten, wobei die maximale Datenrate beim Lesen von 32-bit Werten erreicht wird. In unserem Fall jedoch liefert das Aquisitionsgerät 16-bit breite Daten. Die Geschwindigkeit des Reduktions- und des Expansions-schrittes ist auf Grund der geringen Anzahl an arithmetischen Befehlen durch die Speicherbandbreite limitiert. Dies entspricht einem niedrigen Verhältnis von Rechenoperationen pro Speicherzugriff.

Beide Probleme werden durch die folgende Strategie gelöst: Durch das Zusammenfassen von zwei oder vier, in x-Richtung nebeneinanderliegenden Bildpunkten werden

Speicheroperationen optimiert. Dies führt im Fall von 16-bit Werten zu Speicherzugriffen von nunmehr mindestens 32-bit, wodurch in einem Thread die optimale Speicherbandbreite ausgenutzt und das Verhältnis von Rechenoperationen pro Speicherzugriff gesteigert wird.

### 3.2 Varianzgesteuerter Glättungsfilter

Das varianzgesteuerte Glättungsfilter beschreibt die Tiefpassfilterung eines Bildpunktes entlang der Richtung, die bezüglich einer lokalen Nachbarschaft die minimale Varianz aufweist. Dadurch wird eine kantenerhaltende Rauschreduktion erreicht. Obwohl die Glättungsrichtung und die der minimalen Varianz übereinstimmen, muss die lokale Nachbarschaft des Glättungskerns nicht notwendigerweise der Nachbarschaft der Varianzberechnung entsprechen.

Analog zur Implementierung der Laplace-Hierarchie werden hier auf Grund der lokalen Speicherzugriffe Texturen eingesetzt. Die relativ zum Mittelpunkt definierten Richtungen können aus Effizienzgründen im Constant Memory abgelegt werden, wodurch Mehrfachzugriffe pro Richtung effizient den Constant Cache nutzen.

Eine weitere Optimierung ist das Entrollen von Programmschleifen. Obwohl CUDA hierfür die Compileranweisung `#pragma unroll` zur Verfügung stellt, wurden die vorkommenden verschachtelten Schleifen in unserem Fall nicht vollständig entrollt. Durch den Einsatz von (nicht offiziell unterstützter) Template-Metaprogrammierung konnten wir schleifenähnliche Strukturen nachbilden, die in jedem Fall entrollt werden.

Diese definitive Entrollung der Schleifen ermöglicht einen weiteren Optimierungsschritt: Die relativ zum Mittelpunkt definierten Richtungen der Varianzberechnung können nunmehr im Code als `switch`-Block realisiert werden, wodurch der Zugriff auf den Constant Memory für die Berechnung entfällt. Wählt man für die Glättung und die Varianzberechnung dieselbe Nachbarschaft, so kann die Nutzung des Constant Memory vollständig vermieden werden.

### 3.3 Multiskalarer Rauschglättungsfilter mit Bewegungsdetektion

Ein Bewegungsdetektor kann durch den Vergleich zwischen den verschiedenen Stufen der Laplace-Hierarchie von  $n$  bereits gefilterten Bildern und dem Aktuellen effizient realisiert werden. Durch die Kombination einer Laplace-Pyramide, eines varianzgesteuerten Glättungsfilters und eines Bewegungsdetektors erhält man den multiskalaren Rauschglättungsfilter mit Bewegungsdetektion (siehe Abbildung 1). Dieser ermöglicht eine Rauschreduktion durch das Filtern der oberen Bandpassstufen mit dem varianzgesteuerten Filter, zusätzlich kombiniert mit der zeitlichen Änderung der Bildfolge.

Da der Registerverbrauch dem der optimierten Laplace-Pyramide ähnlich ist, können mehrere Werte auf einmal verarbeitet werden, wodurch ein Leistungszuwachs erzielt wird.

## 4 Evaluierung und Ergebnisse

Aufgrund der Komplexität wurde die Optimierung des varianzgesteuerten Rauschglättungsfilters priorisiert, wodurch allerdings geringe Ungenauigkeiten entste-

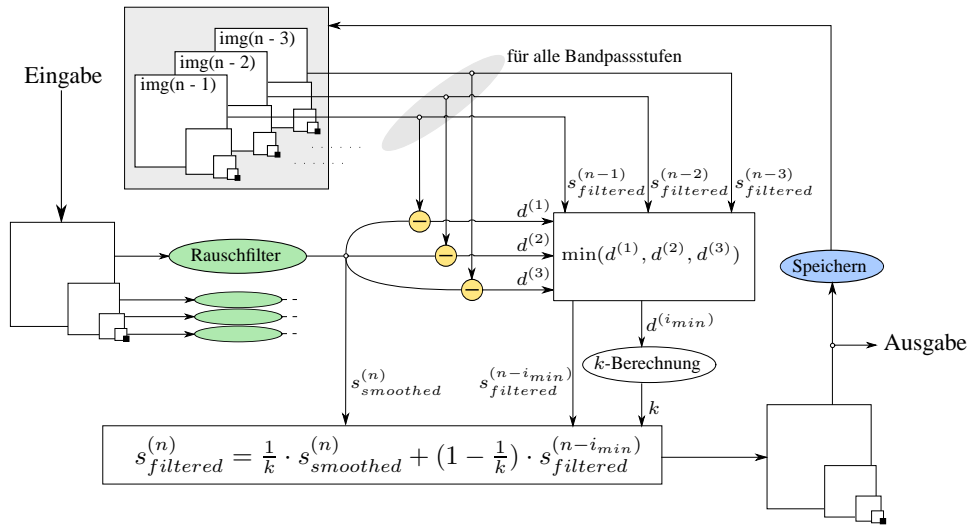


Abb. 1. Verarbeitung der Bandpassbilder mit dem multiskalaren Rauschglättungsfilter

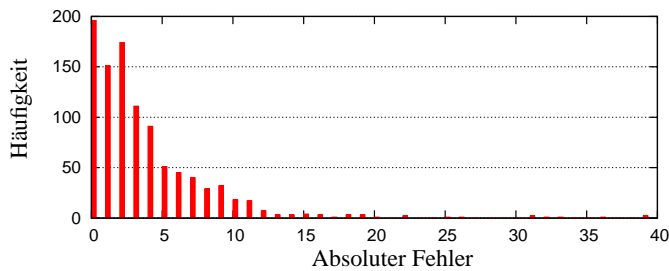


Abb. 2. Häufigkeit absoluter Fehler beim varianzgesteuerten Glättungsfilter, Bild 960x960, 12 Bit Werte

hen (siehe Bild 2). Dabei handelt es sich am häufigsten um geringe Werteabweichungen zwischen der Referenz und der auf der GPU gerechneten Bilder. Insgesamt entsteht dabei im Beispielfeld eine durchschnittliche Abweichung aller Werte von lediglich  $7.2 \cdot 10^{-6}$  für einen Wertebereich von 0 bis 4095. Die fehlerhaften Werte weichen im Schnitt um 1.27% vom Ergebnis der Referenzimplementierung ab.

Der multiskalare Filter hingegen zeigt mit einem durchschnittlichen Fehler von etwa 2.78% stärkere Abweichungen, was wiederum auf die rekursive Struktur des Filters zurückzuführen ist. Zusätzlich entsteht durch die Berechnung mit Fließkommazahlen auf der GPU ein numerischer Unterschied, der jedoch aufgrund der höheren Genauigkeit näher am theoretischen Ergebnis liegt.

Die unterschiedlichen Berechnungszeiten der einzelnen Stufen des multiskalaren Filters sind in Tabelle 1 aufgezeigt. Hierbei verdeutlicht der Vergleich zwischen einer optimierten CPU-Implementierung und der von uns vorgestellten CUDA Implemen-

terung, gemessen auf der für unsere Tests verwendeten Grafikkarte NVIDIA GeForce 8800 GTX, die deutlich schnellere Verarbeitungsgeschwindigkeit. Darüberhinaus liegt die Gesamtleistung mit 7.24 ms in einem Bereich, der einen Echtzeitbetrieb ermöglicht!

## 5 Zusammenfassung

Mit der Einführung von mehrkernbasierten GPUs, wie z.B. NVIDIAs G80 ist es möglich die Anforderungen für Hochleistungsfluoroskopiesysteme in Echtzeit bzgl. der Berechnungsgeschwindigkeit zu erfüllen. Die Verwendung des erweiterten Programmierkonzeptes für Mehrkernarchitekturen (CUDA) und deren 32 Bit Gleitkommaverarbeitung kann die zukünftige Entwicklung von medizinischen Echtzeitanwendungen stark vereinfachen. Allerdings müssen interventionelle Röntgensysteme zusätzlich Netzwerkübertragungen, die Archivierung und verschiedene DICOM-Dienste parallel mit einer minimalen Latenz verarbeiten. Die stabile Hochleistungsverarbeitung eines Bildstroms in Kombination mit einem stark belasteten Prozessor und diesbezüglichen Betriebssystemanforderungen muss daher noch gezeigt werden.

	Glättungs- filter	Kon- struktion Laplace	Rekon- struktion Laplace	Multi- skalarer Filter	Kopieren CPU→ GPU	Kopieren GPU→ CPU	Gesamt
Xeon 2.66 1-Kern	181.59	18.38	15.99	115.35			<b>331.31</b>
Xeon 2.66 2-Kern	91.57	9.88	7.91	58.94			<b>168.30</b>
GeForce 8800 GTX	3.90	0.86	0.48	0.74	0.60	0.66	<b>7.24</b>
Speedup	x46.6	x21.4	x33.3	x155.9			<b>x45.8</b>

**Tabelle 1.** Gemessene Zeiten jedes Filterschrittes in ms bzw. Beschleunigungsfaktor zwischen Singlecore und GPU Implementierung

## Literaturverzeichnis

1. Mueller K, Yagel R. Rapid 3D cone-beam reconstruction with the Algebraic Reconstruction Technique (ART) by utilizing texture mapping graphics hardware. Nuclear Science Symposium, 1998 Conference Record 1998 IEEE 1998;3:1552–1559.
2. Mack M, Hornegger J, Paulus D, Galant A, Böhm S. Echtzeit-Röntgenbildverarbeitung mit Standardhardware. In: Tolxdorff T, Braun J, Handels Heinz, Horsch A, Meinzer HP, editors. Bildverarbeitung für die Medizin 2004. Berlin; 2004. .
3. Scherl H, Keck B, Kowarschik M, Hornegger J. Fast GPU-Based CT Reconstruction using the Common Unified Device Architecture (CUDA). In: Frey EricC, editor. Nuclear Science Symposium, Medical Imaging Conference 2007; 2007. p. 4464–4466.
4. Burt PJ, Adelson EH. The Laplacian Pyramid as a compact image code. IEEE Transactions on Communications 1983;31:532–540.
5. NVIDIA. NVIDIA CUDA - Compute Unified Device Architecture - Programming Guide. NVIDIA Corporation. 2701 San Tomas Expressway, Santa Clara, CA 95050 (USA). version 1.1 ed.; 2007.