

Computations in Graphics Hardware

basierend auf

Using Graphics Cards for Quantized FEM Computations
Virtual 16 Bit Precise Operations on RGBA8 Textures

von Robert Strzodka
und Martin Rumpf

Seminar „Physically-based Methods for 3D Games
and Medical Applications“

Andreas Helbling

Die Autoren



Dipl.-Math. Robert Strzodka

Research Areas

- Hardware sensitive analysis of PDE methods in image processing
- Fast implementations on alternative hardware designs
- Integration of data preprocessing and visualization



Prof. Dr. Martin Rumpf

Working group Numerical Analysis and Scientific Computing

< bersicht

Motivation

Idee

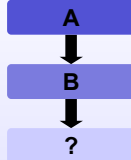
Rechnen

Matrix-Vektor Produkt

Anwendung: Lösen eines Gleichungssystems

Linear Heat Equation

Fire...



Motivation (1)

Vergleich der Architekturen: GPU vs. CPU

	GeForce 4 Ti4600	Pentium 4
Clock	300 MHz	2-3 GHz
Datenbus	128 Bit	64 Bit
Memory	128 MB	Up to 4 GB
Mem clock	650 MHz DDR	533 MHz QDR
Mem latency	2 or 3 clocks	3 clocks
Bandwidth	10.4 GB/s	4.2 GB/s
Caches	Quad Cache, size ?	8kB L1 (data) / 512 kB L2
in Parallel	8 textures / rendering pass	SSE: 4 float ops, MMX: 2 int ops



Hohe Bandbreite
Schnelles RAM
Datenbus ist 128 Bit

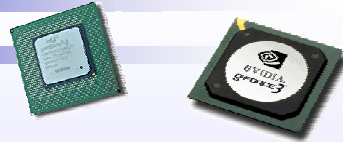
CPU Takt sehr hoch
Viel RAM
Sehr flexibel, kann alles!



GPU Takt tief
Kann nur spezielle Operationen

Effizienz ~ Cacheverhalten

Motivation (2)



GPU

- kommt ohne Caches zurecht
- ausgelegt für grosse Datenmengen
- RAM ist viel grösser als CPU Cache
- spezielle Operationen für grosse Datenmengen
- häufig nicht ausgelastet

CPU

- nur effizient, wenn das Datenvolumen kleiner ist als der Cache
- Operationen, nur auf einzelne/wenige Werte anwendbar
- für Programme optimiert (Branch-Prediction etc..)
- macht alles Andere auch noch

Motivation (3)

Textur-Matrix Analogie

Wie kommt meine GeForce zu Daten ?

• Textur = $[R_1 G_1 B_1 A_1, R_2 G_2 B_2 A_2, R_3 G_3 B_3 A_3, \dots]$

$m \times n$ Texel = $m \times n \times 4$ Farbtensoren aus $[0, 255]$

• Matrix = $[m_{11}, m_{12}, m_{13}, \dots]$

$m \times n$ Zahlen aus $[-\infty, \infty]$



Eine Textur ist eine Matrix [von 4-Vektoren], deren Komponenten nur Integer-Werte aus $[0, 255]$ annehmen können

Idee (1)

OpenGL

- einfache und direkte Nutzung der Grafikkarte
- Hardware-unabhängig
- GLUT: Plattform-unabhängig
- auf vielen Plattformen verfügbar

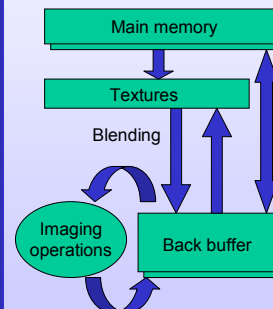
Zentrale Idee: Benutze Blending für die Implementierung von algebraischen Operationen auf Farbwerten

2 Ansätze: Fragment-basierend vs. Texel-basierend

Idee (2)

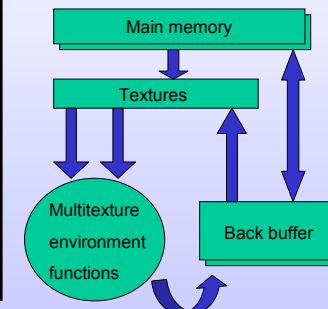
Fragment-basierend:

Benutze Operationen, die direkt auf Pixeln im Framebuffer operieren



Texel-basierend:

Benutze Operationen, die auf Texturen operieren



Rechnen - Einfache Operationen (1)

Fragment-basierend

4

Operation	Formel	OpenGL
addition	$V+W$	BlendFunc
komp. multiplikation	$V \bullet W$	BlendFunc
skalar addieren	$V+bl$	BlendFunc
skalarmultiplikation	aV	BlendFunc
lineare transf.	$aV+bl$	PixelTransfer
funktion	$f(V)$	PixelMap
subtraktion	$V-W$	BlendEquation
lineare transf.	$aV \pm bl$	BlendFunc
funktion	$f(V)$	ColorTable
maximum	$\max(V, W)$	BlendEquation
minimum	$\min(V, W)$	BlendEquation
faltung	$S \ast V$	ConvolutionFilter
vektor norm	$\ V\ _{k=1..n}$	Histogram
color matrix	CV	MatrixMode

$\pm ?$

OpenGL 1.2

V, W : Texturen; a, b : Skalare; I : Identitätsmatrix; C : Matrix

Rechnen - Einfache Operationen (2)

Texel-basierend

Operation	Formel	OpenGL
addition	$V+W$	texture_env_add
komp. multiplikation	$V \bullet W$	standard
lineare transf.	$aV+bl$	texture_env_add
lineare transf.	$aV+bl$	texture_scale_bias
funktion	$f(V)$	texture_color_table
funktion	$f(V_0, V_1, V_2)$	texture_shader
funktion	$f(V_0, V_1, V_2, V_3)$	pixel_texture

V, W : Texturen; a, b : Skalare; I : Identitätsmatrix;

Rechnen - Einfache Operationen (3)

BlendFunc

• überblending

• Fragmente werden gemischt

The `glBlendFunc` function specifies pixel arithmetic.

```
void glBlendFunc(GLenum sfactor, GLenum dfactor);
```

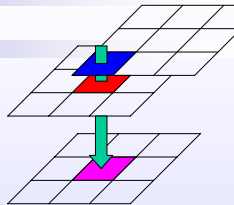
Parameters

sfactor

Specifies how the red, green, blue, and alpha source-blending factors are computed. Nine symbolic constants are accepted: `GL_ZERO`, `GL_ONE`, `GL_DST_COLOR`, `GL_SRC_ALPHA`, `GL_DST_ALPHA`, ...

dfactor

Specifies how the red, green, blue, and alpha destination-blending factors are computed. Eight symbolic constants are accepted: `GL_ZERO`, `GL_ONE`, `GL_SRC_COLOR`, `GL_SRC_ALPHA`, `GL_DST_ALPHA`, ...



Rechnen - Einfache Operationen (4)

BlendFunc Faktoren

• bestimmen resultierendes Fragment

• sind **nicht** beliebig

• sind 1, 0 oder Farbe resp. Alpha einer der Texturen

• Alpha-Werte können beliebig in $[0..255]$ sein

• Farbe/Alpha einer Textur: Faktor **pro Pixel**

```
void glBlendFunc(
    GLenum sfactor,
    GLenum dfactor
);
```

Result = sfactor*Source + dfactor*Destination

Source: Fragmente, die neu in den Framebuffer kommen

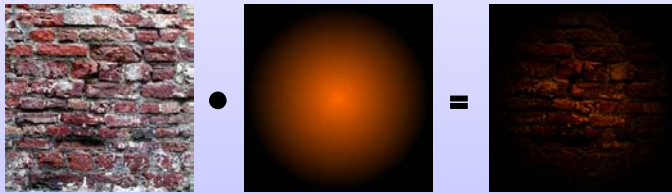
Destination: Fragmente, die bereits im Framebuffer sind

Rechnen - Einfache Operationen (5)

BlendFunc Beispiel: Komponentenweise Multiplikation

```
void display()
{
    glClear(GL_COLOR_BUFFER_BIT);
    glBlendFunc(GL_ONE, GL_ZERO);
    glDrawPixels(256, 256, GL_RGB, GL_UNSIGNED_BYTE, texture);
    glBlendFunc(GL_ZERO, GL_SRC_COLOR);
    glDrawPixels(256, 256, GL_RGB, GL_UNSIGNED_BYTE, spot);
    glFlush();
}
```

```
void glBlendFunc(
    GLenum sfactor,
    GLenum dfactor
);
```



Rechnen - Zahlenformate (1)

Grafikhardware: Intensitäten $\in [0, 255]$ resp $[0, 1]$

→ Wie kann man float-Zahlen mit 8Bit Intensitäten emulieren ?

$$r: [-\rho_0, \rho_1] \rightarrow [0, 1]; \quad r(x) = \frac{(x + \rho_0) \cdot 1}{\rho_0 + \rho_1}$$

$$r^{-1}: [0, 1] \rightarrow [-\rho_0, \rho_1]; \quad r^{-1}(x) = \frac{x}{1}(\rho_0 + \rho_1) - \rho_0$$

513.740483259 ?



• jede Zahl repräsentierbar

• wird nur durch 8 Bit aufgelöst! (Dazu später mehr...)

• ρ_0 und ρ_1 möglichst klein wählen.

• Meistens will man $\rho_0 = \rho_1 = \rho$

Rechnen - Zahlenformate (2)

$$r(x) = \frac{(x + \rho_0) \cdot 1}{\rho_0 + \rho_1}$$

Problem: Wir möchten aus den Zahlen $Z: [-100, 100]$ den Wert $x=50+25=75$ auf der Grafikhardware berechnen. Die Intensität des Resultates sollte also

$$r(75) = \frac{(75+100)}{(100+100)} = \frac{175}{200} = 0.875 \text{ sein.}$$

$$r(50) = \frac{(50+100)}{(100+100)} = \frac{150}{200} = 0.75$$

$$r(25) = \frac{(25+100)}{(100+100)} = \frac{125}{200} = 0.625$$

GPU: $0.75 + 0.625 = 1.375 \rightarrow \text{overflow} \rightarrow \text{😞}$

Trick: $a+b$ rechnet man so: $a+b = 2 \left(\frac{1}{2}r(a) + \frac{1}{2}r(b) - \frac{1}{4} \right)$

$$\rightarrow 2 \cdot \left(\frac{1}{2} \cdot 0.75 + \frac{1}{2} \cdot 0.625 - \frac{1}{4} \right) = 2 \cdot (0.375 + 0.3125 - 0.25)$$

$$= 2 \cdot (0.6875 - 0.25) = 2 \cdot (0.4375) = 0.875 \text{ 😊}$$

Weitere Formeln für Multiplikation, Lineare Transformation, etc...

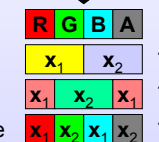
Rechnen - 16 Bit genaue Fixpunktzahlen (1)

Idee:

• Verteile Zahl auf zwei 8Bit Zahlen

• Rechnen wird komplizierter und langsamer

• Implementation z.B. mit PixelShadern



Anforderungen an das 16Bit Format:

• Obermenge der existierenden 8Bit / 9Bit Formate

• Higher 8 Bits = Beste 8Bit Approximation

• Effiziente Implementation (z.B. carry-over)

Rechnen - 16 Bit genaue Fixpunktzahlen (2)

Definition

$$x_1 = \Psi(R,A) = \psi(r,a); \quad R,A \in [0,255]; \quad r,a \in [0,1]$$

$$x_2 = \Psi(G,B) = \psi(g,b); \quad G,B \in [0,255]; \quad g,b \in [0,1]$$

$$\psi(x_{hi}, x_{lo}) = (2x_{hi} - 1) + \frac{1}{128}(x_{lo} - \frac{1}{2})$$

$$\Psi(X_{hi}, X_{lo}) = \frac{1}{255}((2X_{hi} - 255) + \frac{1}{128}(X_{lo} - \frac{1}{128}))$$

Dieses Format erfüllt alle Anforderungen



Rechnen - 16 Bit genaue Fixpunktzahlen (3)

Zusammengefasst:

• H⁺ here Genauigkeit

• F_r spezielle (Teil-)Probleme sinnvoll

• Implementation kompliziert

• Nur wenige Operationen realisierbar

• Neue Operationen nur noch halb so schnell, Multiplikation noch langsamer

Besser warten auf GeForce FX 😊

	$\text{tex0}[RLA, GBE] + (a_1, a_2) \cdot \text{tex1}[RLA, GBE]$
0: RGB	$(\text{tex0}[RGB] - \frac{1}{2}) + (a_1, a_2, a_3) \cdot (\text{tex1}[RGB] - \frac{1}{2}) \rightarrow \text{tex0}[RGB]$
0: A	$\text{tex0}[A] + a_4(\text{tex1}[A] - \frac{1}{2}) \rightarrow \text{tex0}[A]$
1: RGB	$(\text{tex0}[A] < \frac{1}{2})? -((a_1 - 1)/2^8, 0, \frac{1}{2}) : -((a_1 + 1)/2^8, a_2 - \frac{1}{2}) - \frac{1}{2} \text{tex0}[RGB]$
1: A	$\text{tex0}[A] + a_4(\text{tex1}[A] - \frac{1}{2}) \rightarrow \text{tex0}[A]$
2: RGB	$\text{tex0}[RGB] + (0, -1, -1) \cdot \text{tex0}[RGB] \rightarrow \text{tex0}[RGB]$
2: A	$(\text{tex0}[A] - \frac{1}{2}) + a_4(\text{tex1}[A] - \frac{1}{2}) \rightarrow \text{tex0}[A]$
3: RGB	$(\text{tex0}[A] < \frac{1}{2})? -((a_1 - 1)/2^8, 0, \frac{1}{2}) : -((a_1 + 1)/2^8, a_2 - \frac{1}{2}) - \frac{1}{2} \text{tex0}[RGB]$
4: RGB	$\text{tex0}[RGB] + (-1, 0, 0) \cdot \text{tex0}[RGB] \rightarrow \text{tex0}[RGB]$
4: A	$\text{tex0}[A] + (-1) \cdot \text{tex0}[A] \rightarrow \text{tex0}[A]$

Rechnen - Optimierungen (1)

Blending braucht immer mindestens 2 rendering passes

OpenGL Erweiterung `EXT_texture_env_combine` und Texturfunktion `ADD_SIGNED_EXT`

$$a + b = 2 \left(\frac{1}{2} r(a) + \frac{1}{2} r(b) - \frac{1}{4} \right) \text{ in einem rendering pass}$$

Allgemeiner: Erweiterung `NV_register_combiners`

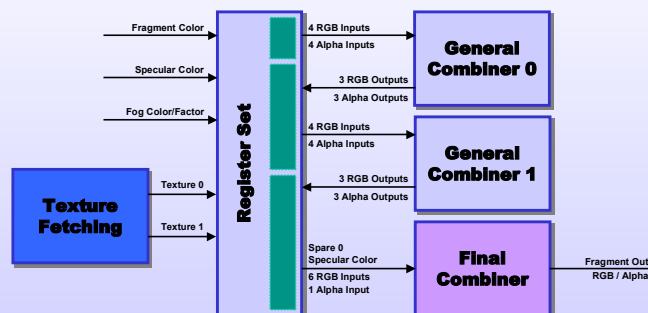
• Addition der Intensitäten implementierbar

• Zwischenresultate der Combiners können in [-1,1] sein

• Addition in einem rendering pass

Rechnen - Optimierungen (2)

Register Combiners



Matrix-Vektor Produkt (1)

Zeile mal Spalte - und zwar komponentenweise

$$A\vec{x} = \vec{y}$$

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} & \dots \\ a_{21} & a_{22} & a_{23} & \dots \\ a_{31} & a_{32} & a_{33} & \dots \\ \dots & \dots & \dots & \dots \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \dots \end{bmatrix} \rightarrow y_i = \sum_j a_{ij} * x_j = \sum_j (\vec{a}_i \cdot \vec{x})_j = \|\vec{a}_i \cdot \vec{x}\|_1$$

komponentenweise Multiplikation von Vektoren

aufsummieren der Komponenten

gleiche Implementation:

A als Menge von Texturen und x als einzelne Textur

komponentenweise Multiplikation mit BlendFunc

1-Norm u.a. mit Histogramm

Anwendung: Lösen eines Gleichungssystems

$Ax=b$

Lösen mit iterativem Algorithmus: $x_{i+1}=F(x_i)$; $x_0=b$;

Jacobi-Iteration

Konjugierte Gradienten

Methoden benutzen nur Operationen, die wir schon gesehen haben

Allgemein: Fast alles ist machbar, nur Frage von Aufwand, Effizienz und Genauigkeit

Linear Heat Equation (1)

zeitabhängige Temperaturverteilung $u(x,t)$

Anfangsverteilung $u_0 = u(t=0)$

Heizungen resp. K, h, k, r per $f(x)$

$$\frac{\partial}{\partial t} u - \Delta u = f$$

Funktion f als diskretisierte Funktion in einer Textur F gespeichert:

$f(x)=f(x_1,x_2)=F[x_1,x_2]$ (look-up table)

Linear Heat Equation (2)

Lösung durch iterative Methode in jedem Zeitschritt

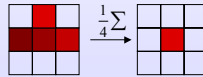
Pseudocode:

```
HeatEquation(u_0, f, lambda)
Load Images representing u_0, f and parameters into Textures
// from now on, everything is done in graphics hardware
for each timestep k
{
    calculate R^k = U^k + lambda F
    Init iterative solver with S^0 = R^k
    for each iteration i
    {
        calculate S^{i+1}
    }
    store U^{k+1} = S^{i+1}
}
```

Fire...

Einfach(st)er Algorithmus für 2D Feuer:

- 1 Addiere an N zufälligen (x,y) in untersten zwei Zeilen einen Wert
- 2 $Img[x,y] = (Img[x,y] + Img[x-1,y] + Img[x+1,y] + Img[x,y+1]) / 4$
- 3 $Img[x,y] > 0 \rightarrow Img[x,y]--;$
- 4 Zeichne Img (Feuer-Palette)
- 5 Goto 1



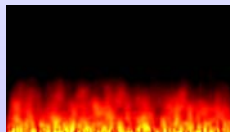
Machbar in Hardware

• Zufallstextur addieren mit BlendFunc

• Summe mit BlendFunc

• $Img[x,y]--$ mit zu kleinen Alpha-Werten

• Palette mit PixelMap



Zusammenfassung

• Neue Perspektive

• Grafikkarte kann große Datenmengen schnell und parallel verarbeiten

• Wichtige Grundoperationen realisierbar

• Grafikkarte als Coprozessor

• Für gewisse Anwendungen geeignet

• Auch komplexere Probleme lösbar (PDEs)

• Relativ kompliziert und beschrankte wissenschaftliche Berechnungen

• Hardware wird immer programmierbarer

Ausblick, Diskussion

Nachteile

• Operationen missen Vorteile der Grafikkarte nutzen → Approximierung von nicht-linearen Funktionen durch lineare

• Intervall größer als $[-1,1]$ → PixelTransfer Funktion, ist aber langsam und eingeschränkt

• Histogramm für Zwischenresultate zu langsam

• Genauigkeit limitiert Anwendungsmöglichkeiten

Ausblick

• PixelShader und RegisterCombiner auf GeForce 3/4

• GeForce FX: 128Bit Farben, d.h. 32 Bit pro Kanal, erweiterte PixelShader