

Implementierung algorithmischer Optimierungen für Volume-Rendering in Hardware

**Entwicklung und Simulation eines Multithreading-
Pipeline-Prozessors zur Visualisierung
dreidimensionaler Datensätze**

Inauguraldissertation
zur Erlangung des akademischen Grades eines Doktors
der Naturwissenschaften
der Universität Mannheim

vorgelegt von
Dipl.-Ing.(FH) Bernd Vettermann
aus Bobenheim-Roxheim

2006

Dekan:	Prof. Dr. Matthias Krause,	Universität Mannheim
Referent:	Prof. Dr. Jürgen Hesser,	Universität Mannheim
Korreferent:	Prof. Dr. Peter Fischer,	Universität Mannheim

Tag der mündlichen Prüfung: 15.11.2006

Danksagung

Herrn Prof. Dr. Reinhard Männer möchte ich danken, dass diese Promotion an seinem Lehrstuhl durchgeführt werden konnte.

Dr. Jürgen Hesser hat durch seine Betreuung und ausgiebigen fachlichen Diskussionen erheblich zum Gelingen dieser Arbeit beigetragen. Für diese Unterstützung bedanke ich mich sehr.

Für die vielen Gespräche und gute Zusammenarbeit möchte ich mich bei den Mitgliedern des Lehrstuhls Informatik V, vor allem der FPGA-Gruppe bedanken.

Der Uni-Mannheim gilt mein Dank für die Durchführung und Finanzierung, des aus dieser Arbeit entstandenen Patentes.

Und nicht zuletzt möchte ich Herrn Prof. Dr. Gerhard Albert und Herrn Prof. Dr. Paul danken, dass ich neben meinen Aufgaben, als Laboringenieur an der Hochschule Mannheim im Institut für Entwurf integrierter Schaltkreise, die Möglichkeit hatte, parallel diese Promotion durchzuführen.

Zusammenfassung

Diese Arbeit beschäftigt sich mit der Beschleunigung von Volumenvisualisierung. Bei der Volumenvisualisierung, auch englisch Volume-Rendering genannt, wird versucht, dreidimensionale Datensätze in einem anschaulichen zweidimensionalen Bild darzustellen. Da dies sehr hohe Ansprüche an Rechenleistung und Speicher stellt, ist es für einen interaktiven Umgang mit den Daten, zum Beispiel bei der Rotation eines gezeigten Objektes, notwendig, zur Beschleunigung spezielle Hardware-Systeme zu entwickeln.

Es werden zuerst die wichtigsten Algorithmen für die Volumenvisualisierung vorgestellt und bereits existierende Volume-Rendering-Systeme erläutert.

Hauptinhalt dieser Arbeit ist die Beschreibung einer neuartigen Architektur für ein Hardware-System zur Echtzeitvisualisierung dreidimensionaler Datensätze mit dem Ray-Casting-Algorithmus. Bei diesem Algorithmus wird von der Bildebene aus, für jeden Bildpunkt ein Sehstrahl durch das Volumen gelegt. Entlang des Strahlverlaufes wird das Volumen in regelmäßigen Abständen abgetastet und für jeden Abtastpunkt eine Reflexion zum Beobachter bestimmt. Auf dem Weg der Reflexionen zum Beobachter wird eine Absorption berücksichtigt und das Restlicht aller Reflexionen aufsummiert. Die Summe entspricht der Helligkeit und Farbe des Bildpunktes.

Für diesen Algorithmus existieren zur Beschleunigung Optimierungstechniken, die zur Erzeugung eines Bildes, nur die unbedingt notwendigen Teile des Volumendatensatzes aus dem Hauptspeicher auslesen. Für die Echtzeitvisualisierung großer Datensätze ist deshalb eine Umsetzung der Optimierungstechniken in Hardware unbedingt notwendig. Durch sie wird allerdings ein wahlfreier Zugriff auf den Speicher notwendig und der Bearbeitungsablauf ist nicht mehr deterministisch, weshalb bisher existierende Hardware-System auf deren Umsetzung verzichten haben.

In dieser Arbeit wird erstmals ein Verfahren vorgestellt, das diese Optimierungstechniken ohne Verluste in Hardware implementiert. Das Verfahren basiert auf der genauen Abstimmung dreier wesentlicher Teile:

1. Einem Pipeline-Prozessor zur parallelen Abarbeitung eines Teilbildes als Multithreading-Architektur. Multithreading bezieht sich hierbei auf den schnellen Wechsel zwischen parallel abzuarbeitenden Sehstrahlen. Hierdurch werden Verzögerungszeiten überbrückt, die bei der Berechnung der Optimierungstechniken und zwischen Adressierung und Datenauslesen des Volumenspeichers entstehen.
2. Einer optimal angepassten Speicherarchitektur, die in den Speicherbausteinen enthaltene Puffer als schnellen Zwischenspeicher (Cache) verwendet und räumlich benachbarte

Volumendaten schneller abrufbar macht.

3. Einer Sortiereinrichtung, die Sehstrahlen bevorzugt bearbeitet, die bereits im Zwischenspeicher liegende Daten verwenden, um zeitraubende Seitenwechsel in den Speichern zu minimieren.

Ziel war es, das System auf programmierbaren Logikbausteinen (FPGA) mit externem Hauptspeicher implementierbar zu machen, um es als Beschleunigerkarte in Standard-PCs verwenden zu können.

Die Effizienz der Architektur wurde über eine C++-Simulation nachgewiesen, wogegen die Implementierbarkeit durch eine Hardware-nahe VHDL-Simulation mit anschließender Synthese überprüft wurde.

Die einzelnen Aspekte der Architektur wurden auf mehreren Konferenzen [39,37,38,41] vorgestellt und in der Zeitschrift IEEE Computer & Graphics [42] veröffentlicht. Weiterhin wurde das Verfahren unter der Nummer WO9960527 beim Deutschen Patent- und Markenamt international angemeldet.

Inhaltsverzeichnis

1	Einleitung.....	11
2	Methoden der Volumenvisualisierung.....	14
2.1	Übersicht.....	14
2.2	Polygon-Rendering.....	16
2.3	Volume-Rendering.....	19
2.3.1	Allgemeines.....	19
2.3.2	Marching-Cube-Verfahren.....	19
2.3.3	Splatting-Verfahren.....	21
2.3.4	Ray-Casting-Verfahren.....	22
2.3.5	Volume-Ray-Tracing-Verfahren.....	24
2.3.6	Vergleich der Volume-Rendering-Verfahren.....	24
3	Optimierungsmöglichkeiten für Volume Rendering.....	26
3.1	Early-Ray-Termination.....	26
3.2	Space-Leaping – Distance-Coding.....	26
3.3	Adaptive-Distance-Sampling-Methode	28
3.4	Datenkohärenz.....	29
3.5	Shear-Warp-Algorithmus.....	30
3.6	Variable Auflösung.....	32
4	Echtzeit-Volume-Rendering Implementierungen.....	33
4.1	Software.....	33
4.2	3D-Texture-Mapping Hardware.....	34
4.3	VIRIM.....	36
4.4	VolumePro – Cube-Familie.....	37
4.5	Vizard II.....	40
4.6	Race II.....	41
5	Ray-Casting-Algorithmus.....	43
5.1	Übersicht.....	43
5.2	Klassifizierung.....	45
5.2.1	Pre-Klassifikation.....	45
5.2.2	Post-Klassifikation.....	46
5.3	Beleuchtungsmodelle, Shading.....	46
5.3.1	Ambiente Reflexion.....	46
5.3.2	Diffuse Reflexion.....	47
5.3.3	Spiegelnde Reflexion.....	48
5.4	Abtastung, Interpolation und Integration.....	49
5.4.1	Absorption.....	52
5.5	Compositing.....	52
6	Hardware-Implementierung der Rendering-Pipeline.....	55
6.1	Pipeline-Architektur.....	55
6.2	Übersicht über die Rendering-Pipeline.....	57
6.3	Adressgenerierung.....	58
6.4	Speicherkontroller.....	60
6.5	Interpolation und Gradientenbestimmung.....	60
6.6	Klassifizierung.....	65
6.7	Shading.....	67
7	Grundprinzip der Architektur für Optimierungen in Hardware.....	71
7.1	Pipelining – Hazard.....	71
7.2	Multithreading zur Steigerung der Auslastung.....	73
7.3	Ausnutzung der Datenkohärenz.....	74
7.3.1	Vergleich der Speicherbausteine.....	75
7.3.2	Single Data Rate SDRAM.....	75

7.3.3	Ausnutzung aktiver SDRAM-Bänke als Zwischenspeicher.....	78
7.4	Implementierung von Multithreading und Sortierfunktion.....	80
7.4.1	Erweiterung des Adressgenerators.....	80
7.4.2	Schaltungsrealisierung der Strahlverwaltung.....	82
7.5	Optimierung der Speicherarchitektur.....	84
7.6	Realisierungsmöglichkeit des Gesamtsystems.....	88
7.6.1	Grundsätzlicher Aufbau.....	88
7.6.2	Kommunikationsablauf.....	89
7.6.3	Abschätzung der Kommunikationszeiten.....	90
8	Ausbaufähigkeit der Architektur.....	93
8.1	Distance-Coding	93
8.2	Perspektive.....	99
8.3	Erweiterungsmöglichkeiten durch die Anwendungs-Software.....	102
8.4	Farbe und zusätzliche Klassifizierung.....	104
8.5	Einbeziehen von Polygondaten.....	105
8.6	Schatten.....	106
8.7	Deformation.....	110
9	Ergebnisse und Diskussion.....	113
9.1	C++-Simulation.....	113
9.1.1	Übersicht.....	113
9.1.2	Untersuchung der Strahlbündelgröße.....	115
9.1.3	Einfluss der Optimierungstechniken auf die Zeilenwechsel.....	116
9.1.4	Unterschiedliche Sortierungskriterien.....	120
9.1.5	Untersuchung der Abhängigkeit verschiedener Größen bei Rotation.....	123
9.1.6	Vergleich unterschiedlicher Speichermodelle.....	128
9.1.7	Unterschiedliche Volumengrößen und Klassifizierung.....	130
9.2	VHDL-Simulation.....	132
9.3	Syntheseergebnisse.....	135
10	Ausblick.....	138
11	Literaturverzeichnis.....	140

1 Einleitung

In nahezu allen wissenschaftlichen Disziplinen kommen mittlerweile bildgebende Verfahren zum Einsatz, die mit Hilfe verschiedener physikalischer Effekte versuchen, die innere Struktur von Objekten erkennbar zu machen. Die bekanntesten sind in der Medizin zu finden. Abbildung 1.1 zeigt beispielsweise einen Computertomograph (CT). Hier werden von allen Seiten Röntgenstrahlen durch die Untersuchungsobjekte geschickt und ihre Abschwächung beim Durchgang gemessen. Daraus können dann Schnittbilder oder auch ein dreidimensionaler Datensatz errechnet werden, dessen Grauwerte in den einzelnen Punkten der Dichte des Objektes entsprechen. Abbildung 1.2 zeigt die CT-Aufnahme eines Herzens.



Abbildung 1.1: Computertomograph[84]

Andere wichtige bildgebende Verfahren sind beispielsweise die Sonographie, die mit Ultraschall arbeitet, oder die Magnet-Resonanz-Tomographie (MRT), die durch Nutzung magnetischer Felder den Kernspin der Atome ermittelt und so Rückschlüsse auf unterschiedliches Gewebe zulässt (Abbildung 1.3)

Weitere Beispiele der Volumenvisualisierung findet man in der Geologie, um Erdbodenstrukturen zu visualisieren, oder in der zerstörungsfreien Materialprüfung, um Einschlüsse oder Risse in Untersuchungsobjekten zu erkennen. Und nicht zuletzt findet man Volumenvisualisierung in der Computergrafik. Da die Volumenvisualisierung aber hohe Ansprüche an Rechnerleistung und

Speicher stellt, steht sie hier in ihrer Verbreitung erst am Anfang.

Ein steigendes Verlangen an einem interaktiven Umgang mit den Volumendaten und an eine bessere Bildqualität, nicht nur im privaten Bereich, ist zu bemerken. So gibt es Ansätze in der Medizin, eine Operation zuvor am Computer durchzuführen. Hierzu muss man allerdings die Daten des Patienten flüssig am Bildschirm bearbeiten und bewegen können. Mit reinen Software-basierten Methoden erreicht man allerdings nur Wiederholraten von 1-10 Bilder pro Sekunde für relativ kleine Datensätze mit der Dimension $256 \times 256 \times 256$ Volumenelementen (256^3).



Abbildung 1.2: Schichtbild eines Herzen mit CT aufgenommen[84]



Abbildung 1.3: MRT-Aufnahme eines Schädels[84]

Die reine Rechenleistung, die für die Visualisierung größerer Bilder und Volumen notwendig ist, könnte durch die immer schneller werdenden Prozessoren und Multiprozessorsysteme zukünftig erreicht werden. Was aber nur sehr viel langsamer zunimmt, ist die Geschwindigkeit, mit der sich die Massenspeicher auslesen lassen. Hier gab es in den letzten Jahren wenig prinzipielle Fortschritte. Zwar wurden neue Speicherschnittstellen entwickelt, die es möglich machen, zwei Datenworte pro Takt auszulesen, indem zur steigenden und fallenden Taktflanke Datenworte abgegriffen werden. Dies ist aber nur beim Lesen großer Blöcke nutzbar, so dass bei einzelnen Leseoperationen doch nur mit einer maximalen Rate von $150 \cdot 10^6$ Datenworte pro Sekunde gelesen werden können.

Welche Auswirkungen dies bei einer Steigerung der Volumengröße hat, wird deutlich, wenn man die Auslesezeiten eines kompletten Datensatzes mit unterschiedlichen Auflösungen betrachtet. Wenn man eine Ausleserate von $150 \cdot 10^6$ Datenworte pro Sekunde annimmt benötigt ein 256^3

großer Datensatz nur 0.11 Sekunden. Damit würde man etwa 10 Bilder pro Sekunde erzeugen können, was noch einigermaßen flüssig auf dem Bildschirm erscheint. Bei einer weiteren Vergrößerung des Datensatzes auf 512^3 , sind es schon 0.89 Sekunden, was etwas mehr als einem Bild pro Sekunde entspricht und bei 1024^3 kann mit 7 Sekunden pro Bild nicht mehr von Echtzeitvisualisierung gesprochen werden.

Diese großen Datenvolumen können somit nicht in näherer Zukunft mit Standard-Computersystemen mit flüssigen Bildwiederholraten visualisiert werden. Glücklicherweise sind aus den Software-Realisierungen Methoden bekannt, die zur Visualisierung nur die unbedingt notwendigen Teile des Volumens auslesen und leere oder verdeckte Bereiche aussparen.

Um somit das Ziel, große Volumina in Echtzeit zu visualisieren mit der aktuell verfügbaren Technik zu erreichen, benötigt man eine speziell für Volumenvisualisierung ausgelegte Beschleunigungs-Hardware. Diese muss einmal die zuvor genannten Methoden der Software-Realisierungen implementieren und andererseits einen speziell angepassten Volumenspeicher besitzen.

Da es für die Volumenvisualisierung viele Methoden gibt, die unterschiedliche Vor- und Nachteile besitzen, bietet es sich an, zur Realisierung der Methoden auf immer komplexer werdende programmierbare Logikbausteine zurückzugreifen. Dadurch können unterschiedliche Methoden und Algorithmen ohne Hardware-Änderungen ausgetauscht werden. Die programmierbaren Logikbausteine sind unter dem Namen FPGA (Field Programmable Gate Array) bekannt. Die größten zur Zeit verfügbaren Bausteine besitzen über 200000 Logikgatter, zusätzliche Multiplizierer und bis zu 1 MByte internem Speicher, die frei programmierbar verwendet werden können.

Bevor näher auf die Implementierung von Volumenvisualisierung auf FPGAs eingegangen wird, soll erst ein Überblick, der verschiedenen Möglichkeiten zur Volumenvisualisierung, gegeben werden.

2 Methoden der Volumenvisualisierung

2.1 Übersicht

Visualisierung bedeutet im Allgemeinen, abstrakte Daten in eine verständliche Form zu bringen und anschaulich darzustellen. In unserem Fall geht es darum Daten, die dreidimensionale angeordnet sind, computergestützt in ein zweidimensionales Bild zur Darstellung, beispielsweise auf einem Monitor, zu bringen.

Im Normalfall handelt es sich bei 3D-Daten um Punkte, Flächen oder Volumeneinheiten im dreidimensionalen Raum, welchen Eigenschaften wie beispielsweise Farbe, Transparenz, Dichte oder andere physikalische Werte zugeordnet werden.

Vor der Visualisierung der Daten müssen diese erzeugt werden, dieser Vorgang wird als *Modellierung* bezeichnet. Der Vorgang, aus solchen Modellen Bilder zu erzeugen, wird *Rendering* oder *Rendern* genannt.

Die zu visualisierenden Daten müssen nicht auf drei Raumdimensionen mit Eigenschaften beschränkt sein, beispielsweise kann eine zeitliche Abhängigkeit der Daten bestehen. Die nachfolgend beschriebenen Verfahren müssen dann allerdings auf jeden 3D-Datensatz zu jedem Zeitpunkt angewandt werden und können so auf den dreidimensionalen Fall reduziert werden. Die Ergebnisbilder könnten danach zu einem Film zusammengesetzt werden.

Etwas schwieriger ist die Darstellung von mehreren Messwerten an einem bestimmten Punkt im Raum; beispielsweise von meteorologischen Daten oder Strömungsmessungen. Hierfür sehr spezialisierte Verfahren, wie *Feature Extraction*, die durch verschiedene grafische Darstellungen mehrere Parameter gleichzeitig darzustellen versuchen, werden in dieser Arbeit nicht berücksichtigt. Die hier gezeigten Verfahren können aber trotzdem als Grundlage zur Visualisierung höherdimensionaler Daten verwendet werden.

Abbildung 2.1 erläutert die Zusammenhänge und begriffliche Einteilung verschiedener Visualisierungsverfahren und deren Ein- und Ausgangsdaten.

Die Darstellungen von 3D-Daten werden in 2 Gruppen eingeteilt:

1. Oberflächendaten, die in Polygondarstellung vorliegen, d.h. Objekte werden nur durch ihre Oberfläche mit ihren Eigenschaften repräsentiert. Die Oberfläche wird durch ein Maschennetz (polygon mesh) aus Polygonen mit Ecken (vertices) und Kanten (edges) aufgespannt.
2. Volumendaten, die einen Raum als geschichtete Bilder beschreiben. Für jedes

Volumenelement (Voxel) liegt normalerweise ein Wert wie Dichte oder Farbe vor. Mathematisch werden sie als dreidimensionale Skalarfelder bezeichnet.

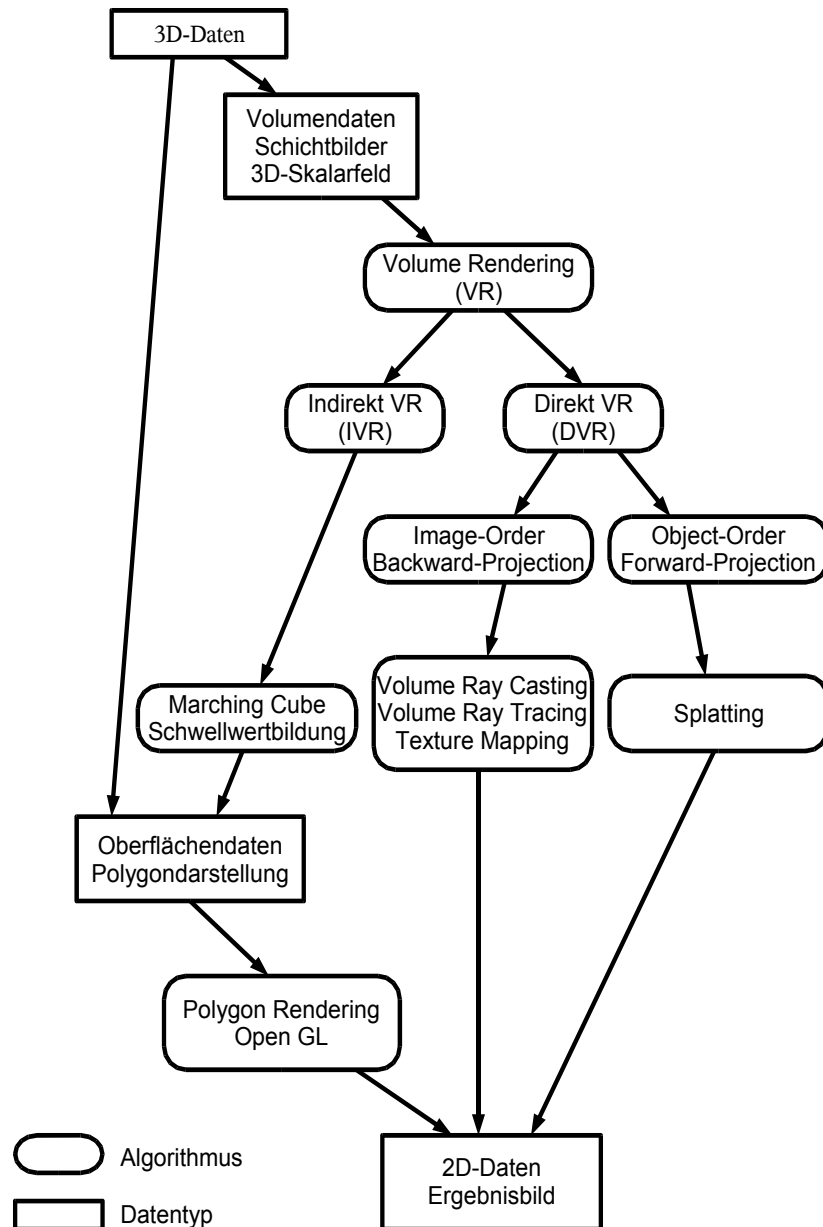


Abbildung 2.1: Zusammenhang zwischen 3D Visualisierungsverfahren und der Datenrepräsentation

Die Verfahren, um diese beiden Repräsentationsformen in zweidimensionale Bilder umzuwandeln, lassen sich danach in *Polygon-Rendering* und *Volume-Rendering* unterteilen.

Um Volumendaten zu visualisieren gibt es zwei *Methoden*:

1. Indirektes Volume-Rendering (IVR), welches über eine Schwellwertbildung versucht dem Volumenobjekt eine eindeutige Oberfläche zu zuordnen und diese dann in eine

Polygondarstellung umzuwandeln. Ein bekanntes Verfahren dieser Gruppe ist als *Marching Cube* [21] bekannt. Die Polygondarstellung wird dann mit Polygon-Rendering weiterverarbeitet.

2. Direktes Volume-Rendering (DVR) erzeugt das Ergebnisbild ohne den Umweg über Polygon-Rendering.

Je nach Richtung der Verarbeitung bei DVR, ob sequentiell über das Ergebnisbild oder sequentiell durch den Speicher, in dem die Volumendaten abgelegt sind, werden die Verfahren unterschieden. Wird für jedes Pixel im Ergebnisbild eine Berechnung durch die Volumendaten gestartet, spricht man von *Image-Order* oder *Backward-Projection*. Den umgekehrten Fall nennt man *Object-Order* oder *Forward-Projection*. Hierbei wird die Berechnung für jedes Volumenelement durchgeführt und sein Beitrag zu einem Bereich des Ergebnisbildes ermittelt.

Bei Object-Order-Verfahren kann optimal in Blöcken und unter Ausnutzung von Speicher-Cache-Verfahren auf die Volumendaten zugegriffen werden, wogegen bei Image-Order-Verfahren algorithmische Optimierungen, die Bestandteil dieser Arbeit sind, leichter und mit höherem Wirkungsgrad implementiert werden können. Es werden dabei wahlfrei einzelne Volumenelemente aus dem Speicher gelesen, was mehr Zeit kostet als in Blöcken zu lesen, dafür muss durch die Optimierungen nur ein kleinerer Teil der Daten gelesen werden.

Im Folgenden werden die einzelnen Verfahren innerhalb dieser Einteilungen genauer beschrieben.

2.2 Polygon-Rendering

Mit Polygon-Rendering-Verfahren werden Oberflächendaten und Bilder zur grafischen Ausgabe verarbeitet. Abbildung 2.2 zeigt ein Beispiel, das absichtlich grob gestaltet wurde, um die Polygone zu erkennbar zu halten. Grafikkarten aller Hersteller unterstützen einen großen Teil dieser Verfahren direkt in Hardware. Als Schnittstelle zur Hardware dienen Programmbibliotheken, die mit standardisierten Aufrufen (API) unterschiedliche Hardware verdecken. Ein solcher Standard ist das OpenGL Software-Interface [13][14].

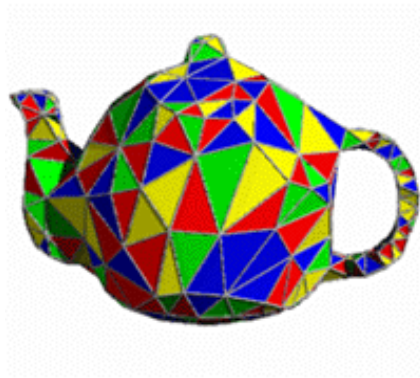


Abbildung 2.2: Teekanne in Polygondarstellung

Abbildung 2.3 zeigt die Repräsentation der Polygone in OpenGL mit Vertices, den Eckpunkten der Polygone [32].

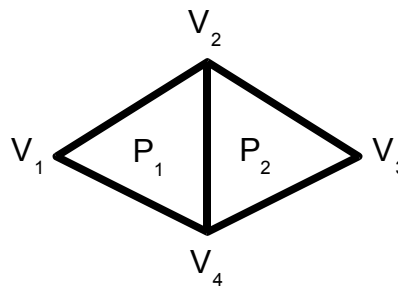


Abbildung 2.3: Polygon-Darstellung in der Oberflächengrafik

Polygone werden als Index in einer Vertex-Liste beschrieben und so in eine für die Grafik-Hardware gut verarbeitbare Form gebracht:

Vertex-Liste:

$$V = (V_1, V_2, V_3, V_4) = ((x_1, y_1, z_1), \dots, (x_4, y_4, z_4)) \quad \text{mit } x_i, y_i, z_i \text{ als Vertexkoordinaten}$$

Polygondefinition:

$$P_n = (V_i, V_j, V_k)$$

Für das Beispiel in Abbildung 2.3 werden die Polygone dann folgendermaßen definiert:

$$P_1 = (1, 2, 4)$$

$$P_2 = (4, 2, 3)$$

Die Vertex-Liste enthält die Koordinaten der Eckpunkte. Zusätzlich kann jeder Vertex noch

Farbdefinitionen, Normalenvektor, Texturkoordinaten oder Materialdefinitionen besitzen.

Die Vertex-Listen und analytisch spezifizierte OpenGL-Primitive (z.B. Kreis) werden mit Hilfe von OpenGL-Kommandos zusammen mit reinen Pixel-Daten an die OpenGL-Pipeline übergeben (Abbildung 2.4). Die Daten können direkt verarbeitet oder in einer Display-Liste für späteren Gebrauch zwischengespeichert werden. Die OpenGL-Primitive werden im Evaluators in reine Vertex-Daten umgewandelt, wodurch dann nur noch Vertex-Daten vorliegen. In der nächsten Stufe werden alle räumlichen Transformationen ausgeführt. Bei der Beleuchtungsberechnung werden den Vertices unter Berücksichtigung der Materialdefinitionen neue Farbwerte zugewiesen. Beim Clipping werden die durch die Vertices beschriebenen Objekte an den Begrenzungsebenen abgeschnitten.

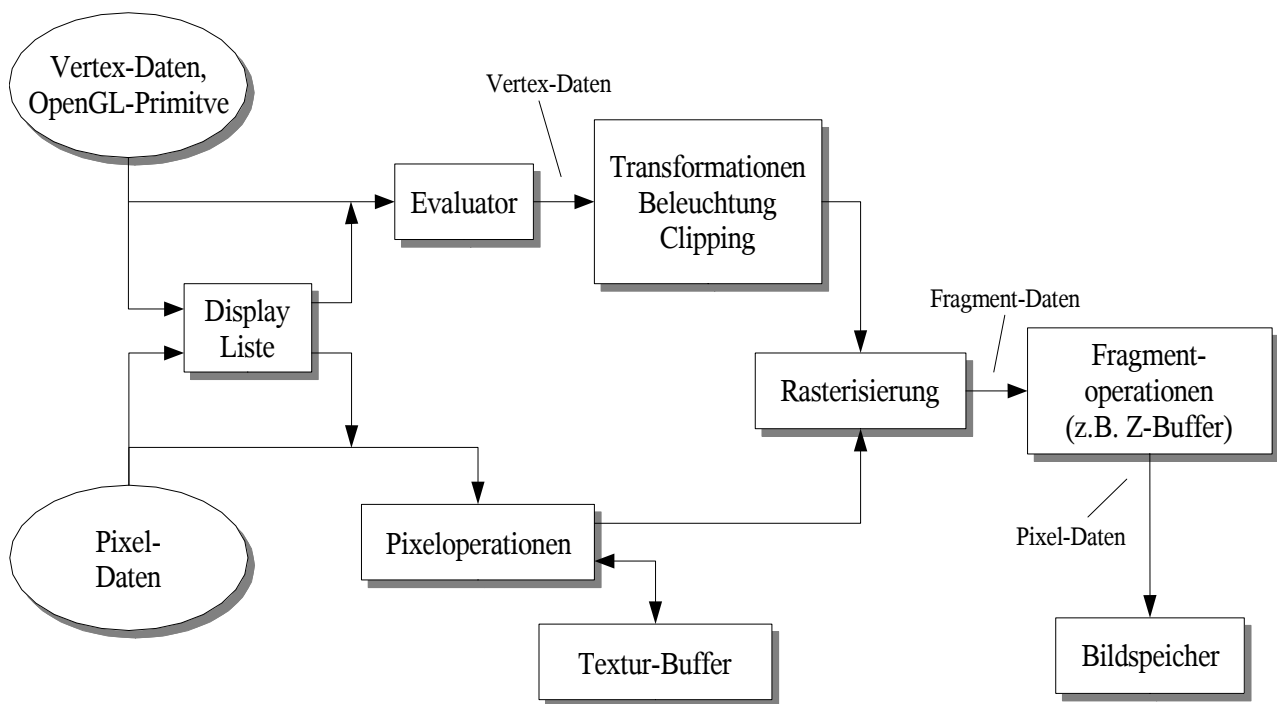


Abbildung 2.4: OpenGL-Pipeline

Für Pixel-Daten und Texturen gibt es spezielle Pixeloperationen. Texturen sind Bilder, die auf die Oberflächen der Vertex-Objekte abgebildet werden. Texturen werden im Textur-Puffer gespeichert und in der Rasterisierung zusammen mit den Vertex-Daten auf einem Raster, das dem des späteren Bildes entspricht, abgebildet. Jedem Rasterpunkt ist nun außer einem Farbwert auch ein Tiefenwert, als Abstand zum Beobachtungspunkt, zugeordnet; man spricht dann von Fragment-Daten. Die wichtigste der darauf folgenden Fragmentoperationen ist die Lösung des Verdeckungsproblems anhand des Tiefenwertes, welches auch Z-Puffer-Test genannt wird. Hierbei wird entschieden welche Bildpunkte im Vordergrund sind und dargestellt werden müssen.

Das Ergebnisbild wird danach in den Bildspeicher geschrieben. Er wird normalerweise als Doppelspeicher verwendet. Während ein Bildspeicher zur aktuellen Anzeige verwendet wird, wird im anderen ein neues Bild gezeichnet und dann umgeschaltet.

2.3 Volume-Rendering

2.3.1 Allgemeines

Volume-Rendering ist der Vorgang, dreidimensionale Datenfelder zu visualisieren.

Ist nicht nur die Oberfläche von Objekten interessant, sondern sollen auch Informationen aus dem Inneren erkennbar werden, reicht Polygon-Rendering zur Visualisierung nicht mehr aus.

Die Information kann oft nur schwer als Anzahl von Polygonen mit Eigenschaften dargestellt werden, vielmehr wird das komplette Objekt in einzelne Schichten mit Volumenelementen zerlegt, die meist einen einzelnen Wert, wie etwa seine physikalische Dichte, als Eigenschaft besitzen. Die Beschreibung der Objekte liegt dann als Volumendatensatz, geschichtete Bilder oder mathematisch ausgedrückt als dreidimensionales Skalarfeld vor.

Beispielsweise werden bei medizinischen Untersuchungen mit Computer-Tomografen (CT) Schichtbilder des Körpers gewonnen. Sie werden dem Arzt in herkömmlicher Weise, als nebeneinander liegende Röntgenbilder, zur Verfügung gestellt. Es erfordert allerdings viel Übung und Erfahrung, diese Einzelbilder in Gedanken wieder räumlich zusammenzusetzen. Soll bei einer Operationsplanung z.B. der ungefährlichste Weg durch Gewebe gesucht werden, ist eine computerunterstützte und interaktive Darstellung der zusammengesetzten Schichtbilder, sehr hilfreich.

Die zur Visualisierung dreidimensionaler Datensätze entwickelten Verfahren werden unter dem Begriff Volume-Rendering zusammengefasst.

Im Folgenden wird das Prinzip einiger wichtiger Verfahren für Volume-Rendering erläutert.

2.3.2 Marching-Cube-Verfahren

Marching-Cube [21] ist ein indirektes Volume-Rendering-Verfahren (IVR) zur Visualisierung von Isoflächen in Volumendaten. Isoflächen sind Flächen, die entstehen, wenn im dreidimensionalen Raum alle Punkte mit gleichen Werten verbunden werden; beispielsweise alle Punkte gleicher Dichte. Die Isoflächen eines vorgegebenen Wertes werden vom Marching-Cube-Algorithmus aus den Volumendaten extrahiert und in Polygonform dargestellt. Diese können dann durch Polygon-

Rendering visualisiert werden. Die Anzahl der darstellbaren Isoflächen ist nur durch den verfügbaren Speicherplatz begrenzt. Die Qualität des Verfahrens hängt sehr stark von den zu visualisierenden Daten ab, da sich nicht alle Objekte durch Isoflächen darstellen lassen.

Die Volumendaten liegen auf einem gleichmäßigen dreidimensionalen Raster. Der Marching-Cube-Algorithmus unterteilt das Raster in 2^3 -Subwürfel und überprüft für jeden Würfel, ob der Wert der acht Volumenelemente über oder unter dem vorgegebenen Schwellwert - dem Isowert – liegen. Zwischen zwei Volumenelementen, die über und unter dem Schwellwert liegen, wird ein Eckpunkt der Isofläche angelegt. Der Eckpunkt kann im einfachsten Fall genau in der Mitte zwischen den Volumenelementen liegen; die Bildqualität lässt sich aber um einiges steigern, wenn die genaue Lage durch lineare Interpolation bestimmt wird. In Abbildung 2.5 bedeuten schwarze Volumenelemente, dass sein Wert unter dem Isowert liegt und weiße liegen darüber. Wurde für alle zwölf Kanten der 2^3 -Subwürfel bestimmt, wo und ob ein Eckpunkt vorhanden ist, spannen diese dann die Polygon-Oberfläche auf.

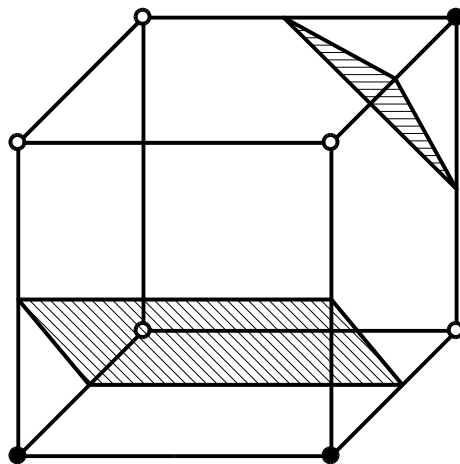


Abbildung 2.5: Bestimmung der Isoflächen im 2^3 Subwürfel

Wenn sich das Objekt durch Isoflächen darstellen lässt, hat der Marching-Cube-Algorithmus den Vorteil, dass die Visualisierung mit Standard-Grafikkarten beschleunigt werden kann.

Probleme entstehen bei der Darstellung von nebelförmigen Gebilden ohne definierte Oberfläche, da der Algorithmus dann keine korrekten, zusammenhängende Isoflächen zuordnen kann. Auch ist keine Information über die Dicke einer Schicht zu gewinnen, da nur deren Oberfläche dargestellt wird [1].

2.3.3 Splatting-Verfahren

Splatting [67] ist ein Direkt-Volume-Rendering-Algorithmus, der als Objekt-Order-Verfahren für jedes Volumenelement den Reflexionsbeitrag zu den Bildpunkten berechnet. Hierzu wird entlang des Sehstrahls von jedem Volumenelement zum Beobachter, der Schnittpunkt mit der Projektionsebene bestimmt. Am Schnittpunkt wird eine dreidimensionale Rekonstruktionsmaske (splat) auf das dort entstehende Bild gelegt. Die einzelnen Werte der Rekonstruktionsmaske werden mit dem Wert des Volumenelementes gewichtet und die einzelnen Beiträge der Maskenelemente zu den, um den Schnittpunkt liegenden, Bildpunkten aufaddiert (Abbildung 2.6).

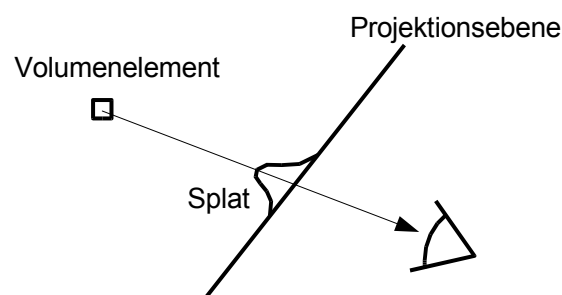


Abbildung 2.6: *Splat: Beitrag eines Volumenelementes zum Bild*

Bildlich kann man sich die 3D-Rekonstruktionsmaske als Schneeball vorstellen, der auf das Bild geworfen wird. Da die Rekonstruktionsmaske räumlich in allen Richtungen symmetrisch ist, kann sie in einer Richtung integriert werden und zu einer zweidimensionalen Maske (Footprint) vereinfacht werden. Dieser Footprint kann vorberechnet werden und in einer Tabelle abgelegt werden, so dass beim Rendern nur noch eine zweidimensionale Operation ausgeführt werden muss. Der Footprint kann als Textur in einer Standard-Grafikkarte hardwarebeschleunigt mit Texture-Mapping auf die jeweiligen Bildpunkte addiert werden (Compositing). Hierbei wird auch berücksichtigt, dass näher an der Bildebene liegende Volumenelemente, gegenüber weiter hinten liegenden, hervorgehoben werden.

Der ursprüngliche Algorithmus lieferte leider keine gute Bildqualität, allerdings wurden in den letzten Jahren einige Verbesserungen vorgeschlagen, die Aliasing-Effekte und Farbfehler vermindern [59,58,68]. So wurde auch die Kantenerkennung durch die Gewichtung mit dem Gradienten möglich und die Qualität durch Post-Klassifikation stark erhöht [66,65] (siehe Kapitel 5.2.2).

2.3.4 Ray-Casting-Verfahren

Ray-Casting [22] ist ein direktes Volume-Rendering und Image-Order Verfahren. Es wurde für diese Arbeit ausgewählt, da man mit ihm eine sehr gute Bildqualität erreichen kann und algorithmische Optimierungsverfahren einfach zu implementieren sind. Hier folgt eine kurze Erläuterung zum Verständnis. Der genauere mathematische Hintergrund wird in Kapitel 5 gegeben.

Im ersten *Abtastungs*-Schritt (Abbildung 2.7) werden von jedem Bildpunkt der Projektionsebene (Ergebnisbild) aus, in Sehrichtung des Beobachters, virtuelle Seh- oder Lichtstrahlen (engl. ray) durch das Volumen „geworfen“ (engl. cast) und entlang der Strahlen in äquidistanten Abständen Abtastpunkte gesetzt.

Da die Positionen der Abtastpunkte normalerweise nicht auf dem Raster der Volumendaten liegen, müssen die Abtastpunkte aus den umliegenden Werten interpoliert werden, um Artefakte, die sich in Form von Würfeln mit der Seitenlänge des Gitterabstandes im Bild wieder finden, zu vermeiden.

Die Sehstrahlen müssen nicht parallel laufen, sondern können, indem sie pyramidenförmig auseinander streben, einen Eindruck von Perspektive vermitteln.

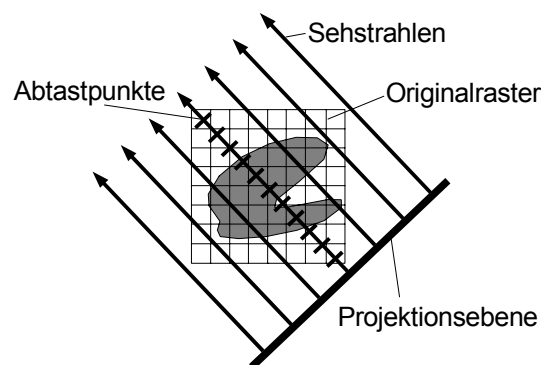


Abbildung 2.7: Abtastung im Ray-Casting-Verfahren

Vor und nach der *Interpolation* kann eine *Klassifizierungsstufe* eingebaut werden. Dadurch können den Originaldaten andere Eigenschaften oder Werte zugewiesen werden. Durch unterschiedliche Gewichtung bestimmter Wertebereiche, können unterschiedliche Details der zu visualisierenden Objekte hervorgehoben werden. Will man beispielsweise bei der CT-Aufnahme eines Kopfes die Schädelknochen sichtbar machen, können niedere Dichtewerte auf Null gesetzt und damit Haut und Muskeln ausgeblendet werden.

Weiterhin gibt es die Möglichkeit, vorab in den Daten, Teile der Volumenobjekte zu kennzeichnen, wie zum Beispiel unterschiedliche Gewebe im Körper, die dann durch die Klassifikation sichtbar oder unsichtbar geschaltet werden können. Man spricht dann von vorsegmentierten Daten.

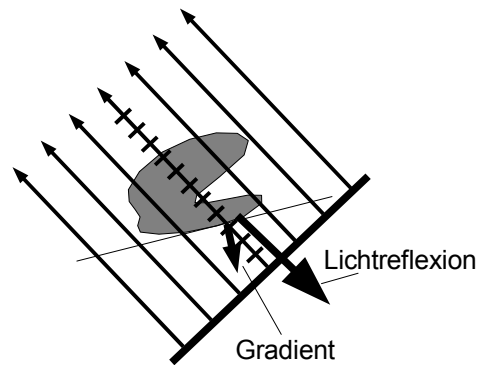


Abbildung 2.8: *Shading im Ray-Casting-Verfahren.*

Nach der Interpolation und Klassifikation wird im nächsten Schritt, dem Shading (Abbildung 2.8), für jeden Abtastpunkt eine Lichtreflexion in Richtung des Beobachters berechnet. Hierzu wird der Gradient an der Stelle des Abtastpunktes bestimmt. Der Betrag des Gradienten gilt als Wahrscheinlichkeitsmaß für das Vorhandensein einer Oberfläche. Mit ihm wird die Reflexion gewichtet, wodurch Kanten der Objekte besser sichtbar werden.

Die Richtung, in die der Gradient zeigt, entspricht der Oberflächennormale und zeigt senkrecht von der Oberfläche weg. Die Oberflächennormale wird für die Berechnung einer diffusen und spiegelnden Reflexion benötigt.

Für das reflektierte Licht wird beim Durchgang der Lichtstrahlen durch die Volumenelemente bis zur Projektionsebene eine Absorption, also eine Abschwächung des Lichtes, berücksichtigt, wodurch Kanten, die näher beim Beobachter liegen, heller erscheinen.

Im letzten Schritt, dem *Compositing* (Abbildung 2.9), werden die Reflexionsbeiträge aller Abtastpunkte entlang eines Sehstrahls aufsummiert. Die Summe entspricht der Helligkeit des Bildpunktes im Ergebnisbild, das auf der Projektionsebene liegt.

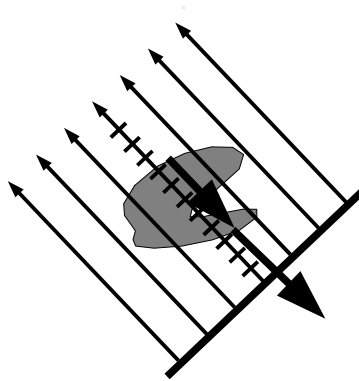


Abbildung 2.9: *Compositing:
Aufsummieren aller Reflexionen*

2.3.5 Volume-Ray-Tracing-Verfahren

Das Ray-Tracing-Verfahren erweitert Ray-Casting um die Darstellung von Schatten. Es wird im Gegensatz zu Ray-Casting nicht nur eine Absorption für die Reflexionen vom Abtastpunkt zum Beobachter berücksichtigt, sondern auch für das eintreffende Licht, wie in Abbildung 2.10 angedeutet. Damit Schatten entstehen, muss es mindestens eine seitliche Lichtquelle geben.

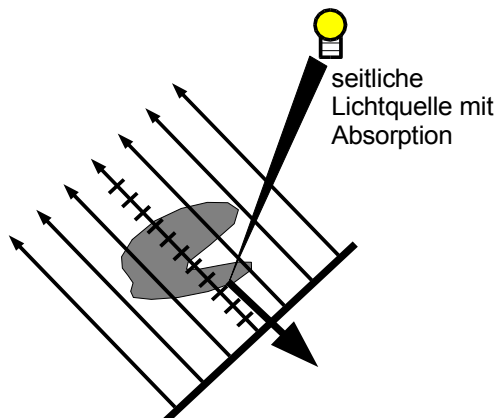


Abbildung 2.10: *Ray-Tracing mit einer zusätzlichen
Lichtquelle*

2.3.6 Vergleich der Volume-Rendering-Verfahren

Im Folgenden sollen die Vor- und Nachteile der gezeigten Volume-Rendering-Verfahren kurz aufgezeigt werden und der geeignetste Algorithmus in Bezug auf allgemeine Verwendbarkeit und Bildqualität ausgewählt werden.

Das Marching-Cube-Verfahren ist eigentlich kein echtes Volume-Rendering-Verfahren, sondern nur

eine Methode aus Volumendaten ein Polygonmodell zu erzeugen. Mit einem Durchgang des Marching-Cube-Verfahrens kann nur eine Oberfläche erzeugt und somit keine semitransparenten Objekte dargestellt werden. Desweiteren können nicht alle Objekte durch exakte Oberflächen dargestellt werden.

Das Splatting-Verfahren erreichte erst durch Forschungsergebnisse in letzter Zeit eine wesentlich bessere Bildqualität. Obwohl für jedes Volumenelement nur 2D-Operationen ausgeführt werden liefern die bekannten Implementierungen allerdings kaum höher Bildwiederholraten, als die von Ray-Casting-Ansätzen. Interessant wird das Splatting-Verfahren, durch die Möglichkeit bessere Interpolationsfilter umzusetzen, als es durch die trilineare Interpolation beim Ray-Casting-Verfahren normalerweise möglich ist.

Das Ray-Casting-Verfahren und mit der Erweiterung um Schatten das Ray-Tracing-Verfahren liefern im Vergleich mit den zuvor genannten Methoden die beste Bildqualität. Sie stellen allerdings auch die größten Anforderungen an die Hardware, will man interaktiv und in Echtzeit Volumenobjekte mit der Ray-Casting-Methode visualisieren.

Für einen 1024^3 Datensatz mit 16 Bit Auflösung werden 2 GByte Speicher benötigt. Sollen diese mit 30 Bildern pro Sekunde visualisiert werden, muss der Speicher mit einer Datenrate von 60GByte pro Sekunde ausgelesen werden können.

Hinzu kommen etwa 300 Milliarden Instruktionen, wenn man nur 10 Instruktionen pro Volumenelement rechnet.

Dies kann momentan keine Standard Hardware erfüllen. Alleine das einmalige Auslesen eines so großen Datensatzes benötigt mit aktuellen Speichern ohne Parallelisierung mehrere Sekunden. Will man dem Ziel, große Datensätze in Echtzeit zu visualisieren, näher kommen, muss man deshalb versuchen, für die Neuberechnung eines Bildes nur die unbedingt notwendigen Teile der Volumendaten auszulesen. Hierzu gibt es bei den rein Software-basierten Systemen einige Methoden zur Datenreduktion und Beschleunigung, die unter dem Begriff algorithmische Optimierungen zusammengefasst werden.

Existierende Software-Systeme können aber höchstens 256^3 -Datensätze in Echtzeit darstellen und erreichen damit etwa 10-20 Bilder pro Sekunde (siehe Kapitel 4.1).

Deshalb müssen spezielle Hardware-Systeme entwickelt werden, die einerseits eine angepasste Speicherarchitektur besitzen und andererseits den Ray-Casting-Algorithmus mit Optimierungen möglichst effektiv unterstützen.

Das folgende Kapitel soll ein Überblick der wichtigsten algorithmischen Optimierungen geben.

3 Optimierungsmöglichkeiten für Volume Rendering

3.1 Early-Ray-Termination

Beim Ray-Casting-Verfahren wird beim Durchgang der Sehstrahlen oder Lichtstrahlen durch das Material eine Absorption berücksichtigt (Siehe auch Kapitel 5.4.1). Fällt die Intensität der Strahlen unter einen Schwellwert, wird die weitere Berechnung des Strahles abgebrochen (Abbildung 3.1), da die nachfolgenden Beiträge zum Bildpunkt der Projektionsebene vernachlässigbar sind. Dieses Verfahren wird Early-Ray-Termination genannt [48].

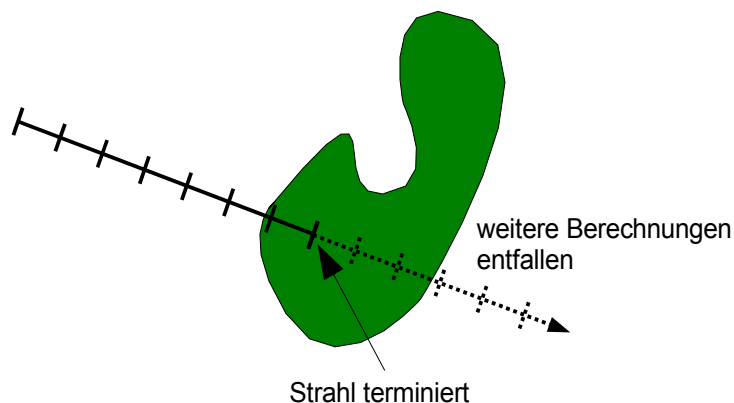


Abbildung 3.1: Early-Ray-Termination

Der Schwellwert kann meist vom Anwender eingestellt werden, um einen Kompromiss zwischen Rechengeschwindigkeit und Genauigkeit der Berechnung zu ermöglichen.

3.2 Space-Leaping – Distance-Coding

Leere Bereiche im Volumen liefern beim Ray-Casting-Verfahren keinen Beitrag zu einem Bildpunkt auf der Projektionsebene, deshalb können diese ohne Änderung des Berechnungsergebnisses übersprungen werden (siehe Abbildung 3.2) [30].

Werden beim Ray-Casting-Verfahren zur Betonung der Kanten, die Reflexionen mit dem Gradienten gewichtet (siehe 2.3.4 und 5.5), kann Space-Leaping sogar auf homogene Bereiche erweitert werden. An diesen Stellen ist der Gradient gleich null, wodurch auch die Reflexionen verschwinden und kein Beitrag zum Bild hinzukommt.

Damit der Algorithmus im voraus die leeren Bereiche bestimmen kann, ist eine Vorberechnung notwendig. Leere Bereiche können einmal durch eine hierarchische Repräsentation der Volumendaten als markierte Blöcke dargestellt werden[48,30] oder als Entfernungsangabe zum

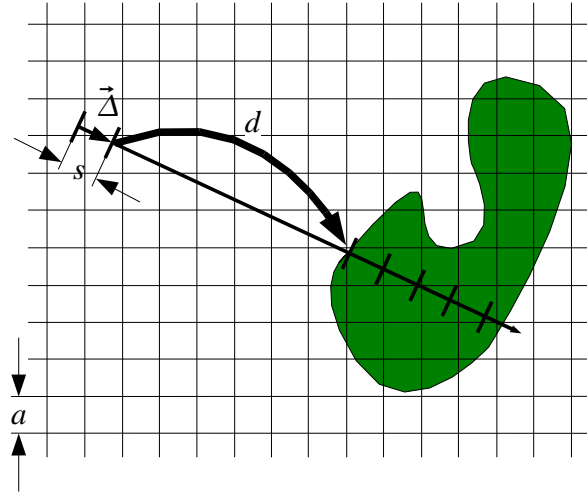


Abbildung 3.2: Space-Leaping

nächsten nicht leeren Volumenelement innerhalb des Volumedatensatzes. Die Berechnung der Entfernungsangabe ist als *Distance-Coding* oder *Distance-Transform* [55,56,16,57] bekannt.

Distance-Coding bestimmt für jedes leere Volumenelement die Sprungdistanz d zum nächsten nicht leeren Volumenelement in der räumlichen Umgebung. Dadurch kann bei der Bildberechnung der nächste Abtastpunkt um diese Distanz in Strahlrichtung weiter geschoben werden, ohne das Ergebnis zu verändern. Da Distance-Coding den Abstand zum nächsten Volumenelement im Raum berechnet, ist die Orientierung des Strahls beliebig.

Der minimale Abstand und die Richtung zweier Abtastpunkte zueinander, wird durch einen Differenzvektor $\vec{\Delta}(X, Y, Z)$ angegeben. Der Betrag des Vektors entspricht der Abtastweite s und ist auf den Gitterabstand a normiert. Die Sprungdistanz wird in Vielfache der Abtastweite umgerechnet, wodurch der nächste Abtastpunkt $\vec{P}_{n+1}(X, Y, Z)$ sehr einfach berechnet werden kann:

$$\vec{P}_{n+1} = \vec{P}_n + \vec{\Delta} d$$

Für die Distanzwerte wird ein getrennter Datensatz angelegt, der genauso strukturiert ist, wie der Volumenspeicher, so dass mit den Koordinaten des Abtastpunktes die zugehörige Sprungdistanz gelesen werden kann. Speicherplatz kann gespart werden, indem die Distanzwerte nicht für jedes Volumenelement, sondern für kleinere Würfel bestimmt werden. Ein Würfel kann dann nur übersprungen werden, wenn er komplett leer ist.

Abbildung 3.3 zeigt die Distanzwerte d in einer Ebene um drei nicht leere Volumenelemente. Die Zahlenwerte geben den minimalen Abstand zu diesen Volumenelementen an. Liegt ein Abtastpunkt zwischen lauter leeren Volumenelementen kann um den zugehörigen Distanzwert weitergesprungen werden.

4	3	3	3	3	3	3	3
4	3	2	2	2	2	2	2
4	3	2	1	1	1	1	2
4	3	2	1	■	■	1	2
4	3	2	1	1	■	1	2
4	3	2	2	1	1	1	2

Abbildung 3.3: Endergebnis von Distance-Coding:
Sprungdistanzen um nicht leere Volumenelemente

3.3 Adaptive-Distance-Sampling-Methode

Während die vorangegangenen Optimierungsmöglichkeiten bei leeren Bereichen (Space-Leaping) und undurchsichtigen Bereichen (Early-Ray-Termination) sehr gut arbeiten, wurde eine Methode unter dem Namen *Adaptive Distance-Sampling* [41] veröffentlicht. Sie ist zwar aufwändiger, kann aber auch semitransparente Bereiche mit einbeziehen.

Die Sprungweiten in *Adaptive-Distance-Sampling*, werden so berechnet, dass durch Auslassen von Abtastpunkten ein vorgegebener Fehler bei der Integration nicht überschritten wird. Dabei wird die Opazitätsänderung, zwischen den Abtastpunkten, linear interpoliert und mit den exakt berechneten Werten verglichen. Abbildung 3.4 verdeutlicht dies in einer Dimension. Die schraffierte Fläche stellt die integrierte Funktion über jedes Abtastpunkt dar, während bei der grauen Fläche die beiden mittleren Abtastpunkte ausgelassen wurden. Wird bei der Differenz der beiden Flächen der vorgegebene maximale Fehler nicht überschritten, dürfen die beiden Abtastpunkte übersprungen werden.

Damit die Sprungweiten richtungsunabhängig berechnet werden, muss der Fehler räumlich um jedes Volumenelement eingehalten werden. Bildlich wird um jeden Gitterpunkt ein Kugel immer weiter vergrößert, bis sich gerade keine Gitterpunkte in der Kugel befinden, für die der Fehler zu groß wird. Der Radius der Kugel entspricht dann der richtungsunabhängigen Sprungweite für den Rendering-Algorithmus.

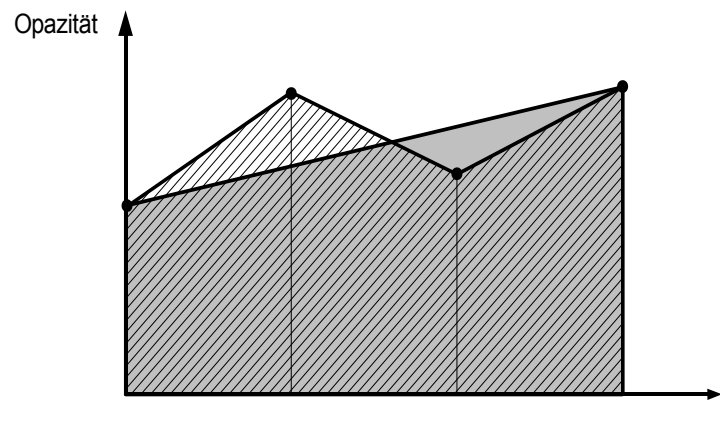


Abbildung 3.4: Fehlerbestimmung bei der Adaptive Distance-Sampling-Methode

Bei der tatsächlichen Berechnung wird die maximale Sprungdistanz beschränkt, so dass nur eine kleine Umgebung um jedes Volumenelement berücksichtigt werden muss. Eine maximale Sprungdistanz von sieben zeigte sich in den Experimenten als ausreichend. Weiterhin kann die Rechnung vereinfacht werden indem der Fehler nur in bestimmten Richtungen überprüft wird. Denkt man sich das Volumenelement im Mittelpunkt eines Würfels, so reichen in den meisten Fällen die drei Hauptachsen und die vier Diagonalen im Würfel zur Überprüfung aus.

Mit der *Adaptive-Distance-Sampling-Methode* können vor allem homogene und nur leicht fluktuierende halbtransparente Regionen übersprungen werden. Die Sprungweiten für leere Bereiche können weiterhin mit dem einfacheren Distance-Coding berechnet werden.

Nachteilig wirkt sich aus, wenn Farbinformation mit verarbeitet wird, dass dann beide, der Fehler der Opazitätsänderung und der der Farbänderung, berücksichtigt werden müssen. Ebenso geht der Geschwindigkeitsvorteil bei verrauschten Datensätzen stark zurück.

3.4 Datenkohärenz

Bei den Volumendaten handelt es sich normalerweise um große Datenmengen, die in dynamischen Speichern abgelegt werden müssen. Diese Speicher sind in Blöcke unterteilt, wobei Zugriffe innerhalb der Blöcke wesentlich schneller von statten gehen, als über verschiedene Blöcke hinweg. Die Blöcke müssen erst bereitgestellt werden, was Zeit kostet und neu zu ladende Blöcke können bereits geladene verdrängen. Für ein optimiertes Rendering-System muss man einen ständigen Wechsel zwischen sich gegenseitig verdrängenden Blöcken vermeiden.

Unter Datenkohärenz versteht man, Daten, die vom Algorithmus kurz nacheinander benötigt werden auch soweit als möglich im gleichen Speicherblock abzulegen. In unserem Fall werden

Daten, die räumlich nah beieinander liegen, bevorzugt verwendet. Die Daten dürfen deshalb nicht ebenen- oder gar zeilenweise in die Speicherblöcke geladen werden, sondern müssen als möglichst großer Würfel im Speicherblock liegen. Dadurch wird eine Richtungsunabhängigkeit bei der Abarbeitung der Daten gewährleistet.

Würde man in einer höheren Programmiersprache einfach ein dreidimensionales Feld deklarieren und die Volumendaten darin abspeichern, würde dies einer zeilenweisen Anordnung in den Blöcken entsprechen. Bei einer Auslesereihenfolge entlang der Zeilen in X-Richtung würde dies funktionieren, aber in Y- wie in Z-Richtung hätte jeder Lesezugriff ein Neuladen von Blöcken zur Folge.

3.5 Shear-Warp-Algorithmus

Der Shear-Warp-Algorithmus ist ein Object-Order-Algorithmus, da er die Volumendaten der Reihenfolge nach abarbeitet. Der Shading- und Compositing-Schritt ist allerdings wieder identisch mit dem Ray-Casting-Algorithmus.

Das Shear-Warp-Verfahren nutzt als Object-Order-Algorithmus die Datenkohärenz optimal aus. Der Datensatz muss nur einmal in der abgespeicherten Reihenfolge ausgelesen werden.

Das Ergebnisbild wird über ein Zwischenbild (intermediate Image) errechnet, das parallel zu der Seite des Datenwürfels liegt, die dem Beobachter am meisten zugewandt ist (siehe Abbildung 3.5). Statt eine 3D-Transformation durchzuführen, wird das Datenvolumen in Schichten parallel zum Zwischenbild zerlegt, die gegeneinander "gescher" (shear) werden.

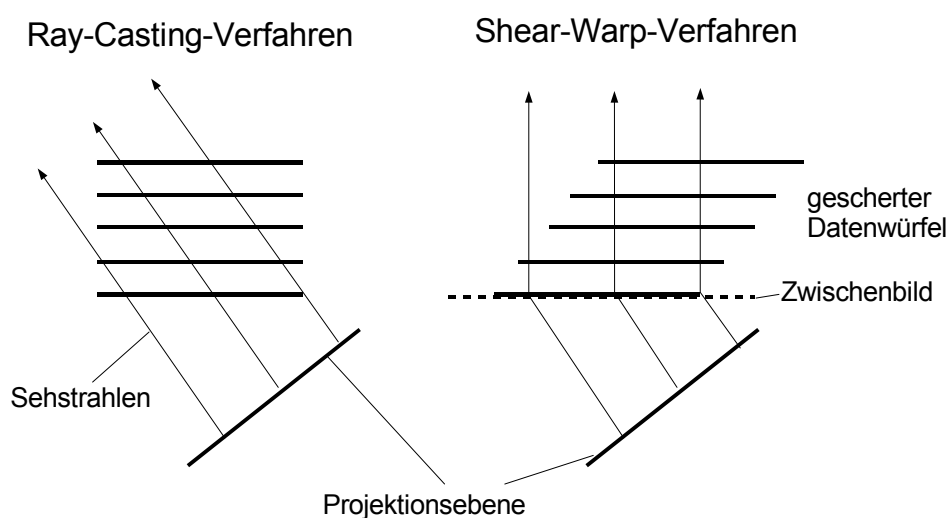


Abbildung 3.5: Projektion beim Shear-Warp-Verfahren

Innerhalb dieser Schichten werden durch 2D-Interpolationen die Abtastpunkte entlang von zwischen den Volumenelementen liegenden Linien, den Voxel-Scanlines, bestimmt. Die übereinander liegenden Voxel-Scanlines werden durch Shading und Compositing auf das Zwischenbild projiziert. Das Zwischenbild muss noch entzerrt werden (warp), was wiederum durch 2D-interpoliertes Abtasten geschieht und das Ergebnisbild liefert.

Abbildung 3.6 zeigt die drei wesentlichen Schritte des Shear-Warp-Algorithmus:

1. Das Scheren des Volumens und Abtasten der Voxel-Scanlines
2. Die Projizierung und Compositing auf den Image-Scanlines des Zwischenbildes
3. Entzerren durch Neuabtastung des Zwischenbildes

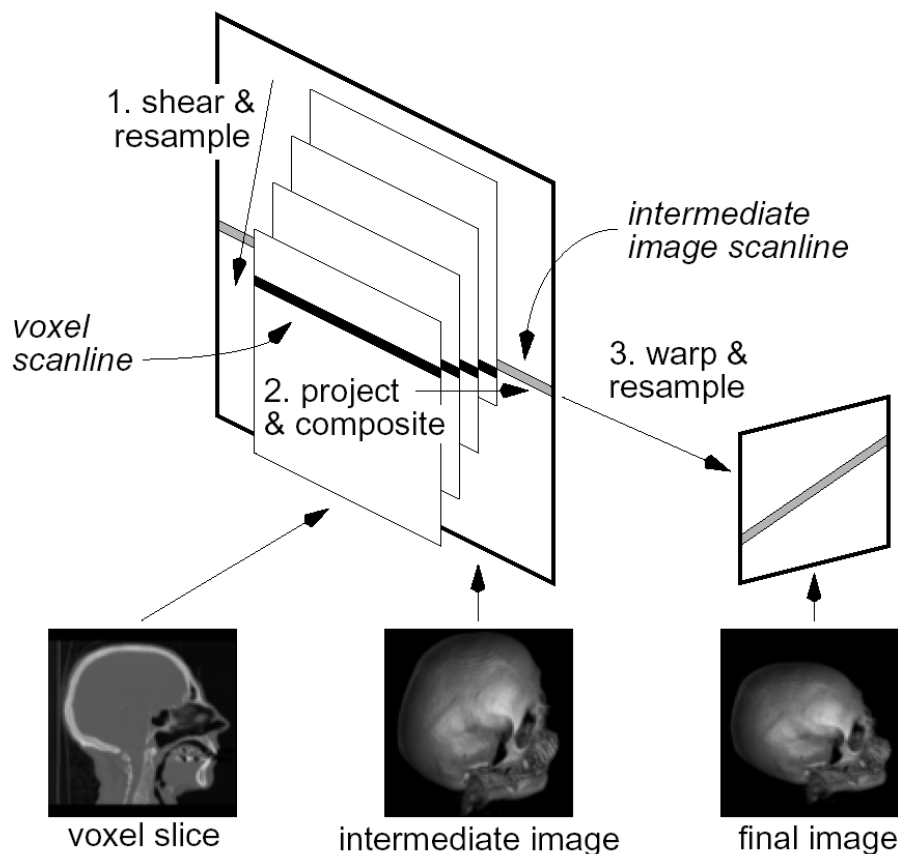


Abbildung 3.6: Die drei Schritte des Shear-Warp-Algorithmus aus [23]

Als algorithmische Optimierungen gibt es bei Shear-Warp das Run-Length-Encoding, das ähnlich dem Space-Leaping, Sprunginformationen zu den nächsten nicht leeren Volumenelementen ablegt. Im Gegensatz zu Space-Leaping gilt die Sprunginformation aber nur innerhalb einer Voxel-

Scanline. Durch Weglassen der leeren Volumenelemente und Hinzufügen der Sprunginformation kann die Datenmenge etwas reduziert werden, obwohl drei getrennte Datensätze für jede Beobachterraichtung in einer Vorberechnung erzeugt werden müssen. Run-Length-Encoding kann innerhalb der Voxel-Scanline und innerhalb der Image-Scanline angewendet werden. In [23] wird das Verfahren innerhalb der Image-Scanline zwar mit Early-Ray-Termination verglichen, es kann aber nie seine Effizienz erreichen, da es nur auf dem zweidimensionalen Zwischenbild angewendet wird, während Early-Ray-Termination im Volumenbereich arbeitet.

Nachteilig auf die Bildqualität wirken sich einmal die zwei Abtastungsschritte aus, die zu Verwaschungen führen, zum anderen findet bei der Abtastung nur eine lineare Interpolation innerhalb der Ebene, in der die Voxel-Scanlines liegen, statt. In orthogonaler Richtung wird die nächstliegende Ebene (Nearest-Neighbour-Interpolation) verwendet, wodurch Artefakte erzeugt werden. Bei der Rotation des Volumenobjektes ist es daher meist erkennbar, wenn das Zwischenbild auf eine andere Seite des Datenwürfels gelegt wird und es kommt zu einem Flackereffekt.

Wegen dieser Nachteile und der Schwierigkeit Early-Termination in einem Object-Order-Verfahren effizient zu implementieren wurde auf die Umsetzung von Shear-Warp verzichtet.

3.6 Variable Auflösung

Variable Auflösung ist eine einfache Methode bei interaktiven Anwendungen die Rendering-Zeit zu verkürzen. Während der Benutzer ein visualisiertes Objekt bewegt, wird das Volumen nur grob, zum Beispiel mit halber Abtastrate, abgetastet. Dadurch wird nur noch ein Achtel der Abtastpunkte erzeugt, wodurch die Bildwiederholrate erhöht wird. Die geringere Auflösung fällt, während das Objekt bewegt wird, kaum auf, aber die Bewegung wird flüssiger. Erst wenn das Objekt nicht mehr bewegt wird, wird es mit voller Auflösung berechnet.

4 Echtzeit-Volume-Rendering Implementierungen

4.1 Software

Es gibt nur wenig Versuche ein Echtzeitsystem für Volume-Rendering in Software auf einer Standard-Hardware zu implementieren. Eines der bekanntesten ist das von Lacroute [24], der den Shear-Warp-Algorithmus [23] (Kapitel 3.5) auf einer 16 Prozessor, Shared Memory, SGI-Challenge mit 13 Hz Bildwiederholfrequenz für einen 128^3 Datenwürfel und ein 256^2 Bild implementiert hat. Das System erreicht etwa 1 Hz Bildwiederholfrequenz für 256^3 Datensätze, allerdings sind dafür einige Vorberechnungen für Run-Length-Decoding notwendig und die Darstellung ist auf Parallelprojektion beschränkt. Eine der letzten Implementierungen des Shear-Warp-Algorithmus erreicht auf einem 3GHz-Pentium IV PC für Datensätze der Größe $256 \times 265 \times 128$ etwa 10 Hz Bildwiederholrate[25].

Das UltraVis System [75,76,77] erreicht auf einem HP Kayak XU PC mit zwei 500MHz Pentium III-Prozessoren und 1 GByte Hauptspeicher etwa 2-10 Hz Bildwiederholfrequenz. Die Volumen sind auf 256^3 begrenzt. Implementiert wurde ein Ray-Casting-Algorithmus mit trilinearer Interpolation, Space-Leaping und Early-Ray-Termination. Die Geschwindigkeit wird erreicht, indem wichtige Teile des Algorithmus, wie Interpolation und Shading, mit Assembler programmiert wurden, um speziell die SIMD^I-Erweiterungen der Intel x86-Architektur auszunutzen. Der begrenzende Faktor des Systems liegt im Datendurchsatz des Hauptspeichers. Deshalb werden die Volumendaten und Parameter so auf den Hauptspeicher verteilt, dass oft verwendete Daten nie aus dem Level 1-Cache des Prozessors verdrängt werden. Als Folge davon muss eine Datenstruktur angelegt werden, die vier mal größerer als der Volumendatensatz ist.

Bildwiederholraten zwischen 1-3 Hz für 256^3 -Volumen werden auch in [78] von Roettger präsentiert. Auch hier handelt es sich um einen Ray-Casting-Algorithmus mit Space-Leaping und Early-Ray-Termination. Erreicht wird dies durch Auslagerung großer Teile des Algorithmus auf den Grafikprozessor einer ATI Radeon 9700 Grafikkarte. Die Programme wurden in der Assembler ähnlichen Sprache *Pixel Shader 2.0* geschrieben. Durch die Implementierung von Pre-Integration konnte auch eine gut Bildqualität erreicht werden.

In [62] wurde das Splatting-Verfahren auf einer NVIDIA-GeForce4-Grafikkarte implementiert, wobei aber nur für kleiner Volumen (128^3) interaktive Bildwiederholraten erreicht wurden.

Obwohl die vorgenannten Verfahren auf Grafikkarten implementiert sind, unterscheiden sie sich

^I SIMD: Single Instruction Multiple Data, Befehlssatz mit speziellen Registern zur Verarbeitung großer Integer-Datenmengen (MMX: Multi Media Extension) und Floatingpoint-Datenmengen (SSE: Streaming SIMD Extension)

deutlich von den im folgenden Kapitel beschriebenen 3D-Texture-Mapping-Verfahren. Dort handelt es sich um *Slice-Based*, also schichtbasierte Verfahren, während hier die Verfahren strahlweise auf dem Grafikprozessor (GPU^I) implementiert sind. Ein Vergleich beider Verfahren konnte zeigen, dass strahlbasierte Verfahren mit höherer Genauigkeit und somit geringeren Artefakten implementiert werden können als schichtbasierte [61]. Hier stellte sich heraus, dass trotz besserer Umsetzung von Early-Ray-Termination und Space-Leaping, die strahlbasierten Verfahren nur etwa die halbe Bildwiederholrate erreichen, da Sprung- und Schleifenunterstützung auf den Grafikprozessoren zwar seit kurzem vorhanden, aber noch unzureichend ist. Ältere Implementierungen mussten deshalb mit mehrfachen Durchgängen arbeiten[78,79,80].

4.2 3D-Texture-Mapping Hardware

Mit 3D-Texture-Mapping wird unter Ausnutzung des Textur-Puffers von Grafikkarten ein vereinfachter Ray-Casting-Algorithmus implementiert. Der Textur-Puffer wird normalerweise dazu verwendet, bei Polygone-Rendering ein Bild (Textur) auf Polygone zu legen (vgl. Polygon-Rendering Kapitel 2.2 und [13]). Zusätzlich bieten die meisten Grafikkarten auch die Möglichkeit übereinander liegende Texturen miteinander zu verrechnen (Blending). Die Texturen werden hierzu entsprechen der Beobachterraichtung neu abgetastet und Zwischenwerte interpoliert. Das fertige, zusammengesetzte Bild kann dann aus dem Frame-Puffer der Grafikkarte ausgelesen werden.

Genau diese Funktionen macht man sich beim 3D-Texturemapping zunutze. Das gesamte Datenvolumen muss hierzu im Textur-Puffer liegen und wird wie beim Ray-Casting neu abgetastet, allerdings sogar schichtweise und nicht strahlweise. Zwischenwerte werden ebenfalls trilinear interpoliert und die Schichten wie beim Compositing überlagert, wie Abbildung 4.1 verdeutlicht. Allerdings und das ist der größte Nachteil, ist direkt in Hardwareunterstützung kein Shading möglich.

Die ersten Implementierungen erreichen zwar bis zu 10 Bilder pro Sekunde für ein 512^2 Bild, verzichten aber vollständig auf Shading und verwenden direkt die Dichtewerte des Volumens. Entsprechend schlecht ist dadurch die Bildqualität [70,71]. Durch neue Funktionen auf modernen Grafikkarten, wie Early-Z-Culling lassen sich Beschleunigungsverfahren wie Early-Ray-Termination und Space-Leaping implementieren, wodurch die Bildwiederholrate für Volumen der Größenordnung 256^3 auf rund 30 Hz gesteigert werden konnte[63,60].

I GPU: Graphical Processing Unit

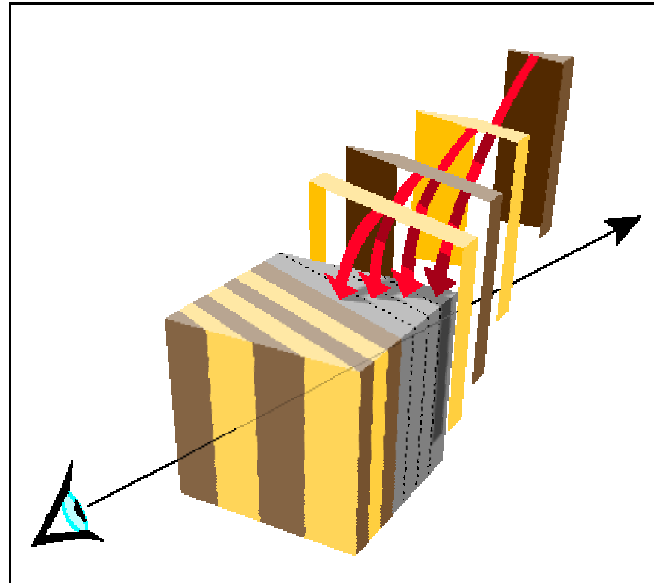


Abbildung 4.1: Schichten bei Volume-Texture-Mapping

Es wurden verschiedene Möglichkeiten veröffentlicht, die Bildqualität weiter zu verbessern. Einmal kann im gesamten Volumen als Vorberechnung Shading durchgeführt werden. Die Shading-Ergebnisse kommen dann in den Textur-Puffer. Die Abtastung und das Compositing wird in der Grafikkarte ausgeführt. Da das Shading aber richtungsabhängig ist, muss für jedes Bild das Volumen neu berechnet und in den Textur-Puffer geladen werden. Pro Bild wurden Berechnungszeiten von 13.4 Sekunden angegeben [72].

Eine weitere Möglichkeit ist, nur die Abtastung mit der Interpolation von der Hardware berechnen zu lassen und das Shading in Software zu berechnen. Hierzu können nicht nur die Dichtewerte interpoliert werden, sondern auch die drei Komponenten vorberechneter Gradienten, die parallel in identischer Weise in den 4 Kanälen $RGB\alpha$ der Grafikkarte abgelegt werden können. Das Compositing kann wahlweise in Software oder im Frame-Puffer durchgeführt werden, indem die Ergebnisse des Shading wieder in den Textur-Puffer zurückübertragen werden. Da beim Software-Compositing Early-Ray-Termination angewendet werden kann, ist es datenabhängig, welche Version schneller ist. Für größere Bilder (512^2) benötigen aber beide Methoden mehrere Sekunden pro Bild [69].

Der größte Vorteil von Volume-Texture-Mapping liegt in der gemeinsamen Verarbeitung von Polygon und Volumendaten in der gleichen Hardware. Für interaktive Anwendungen mit größeren Volumen macht allerdings die fehlende Hardwareunterstützung für Shading und Klassifizierung diese Methode unbrauchbar. Allerdings wurden bereits Erweiterungen für Standard-Grafikkarten vorgeschlagen, um diese Lücke zu schließen [73].

4.3 VIRIM

VIRIM ist bisher das einzigste System, dass Ray-Tracing in Hardware implementiert [3, 2, 4]. Es ist in der Lage Volumen der Größe $256 \times 256 \times 128$, mit 10 Bildern pro Sekunde, mit Perspektive und Schatten zu visualisieren. Es ist modular aufgebaut, wobei immer eine Geometrieeinheit und eine Rendering-Einheit zusammengehören. Alle Module werden über einen VME-Bus von einem Host-Rechner initialisiert und ausgelesen (Abbildung 4.2).

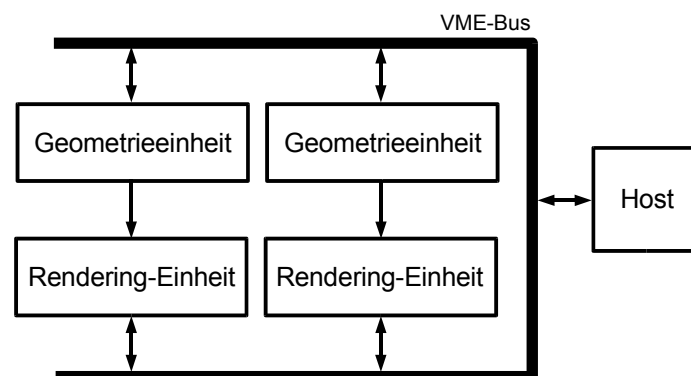


Abbildung 4.2: Komponenten des VIRIM-Systems

Jede Geometrieeinheit enthält jeweils den kompletten Volumendatensatz und gibt schichtweise interpolierte Abtastpunkte mit zugehörigem X,Y-Gradienten an die Rendering-Einheit weiter. Auf den Z-Gradienten wurde wegen des zu großen Aufwandes für die Zwischenspeicherung der Schichten verzichtet. Durch trapezförmige Abtastung der Schichten ist auch eine perspektivische Darstellung möglich.

Auf der Rendering-Einheit werden mit digitalen Signalprozessoren Absorption, Shading und Compositing berechnet. Die Beschleunigung und Skalierbarkeit des Systems wird erreicht, indem pro Rendering-Einheit acht Prozessoren jeweils eine Schicht des Volumens und somit eine Zeile des Ergebnisbildes berechnen. Damit die Prozessoren unabhängig voneinander arbeiten können, gilt als Einschränkung, dass die Strahlen der Lichtquellen (0° und 45°) in der Ebene der Schichten liegen. Die Objekte sind also immer von vorne und von der Seite beleuchtet.

Da es sich um eine feste, mit diskreten Bauteilen aufgebauten, Geometrieeinheit handelt, können ohne neue Hardware-Entwicklungen keine algorithmischen Optimierungen wie Space-Leaping oder Early-Ray-Termination integriert werden. Wogegen durch die Programmierbarkeit der Rendering-Einheit die Implementierung verschiedener Rendering-Verfahren möglich ist.

4.4 VolumePro – Cube-Familie

VolumePro ist eine PCI-Karte, die zusammen mit einer Grafikkarte vertrieben wird und in Standard-PCs eingebaut werden kann. Mit diesem System können 512^3 Volumendaten mit 30 Bildern pro Sekunde visualisiert werden. VolumePro wird von TeraRecon Inc. [11] vertrieben und basiert auf der Cube-4 [8] und EM-Cube-Architektur [9], die in SUNY Stony Brook entwickelt wurden. Kernstück des VolumePro ist der vg500 Rendering-Chip, der 4 parallele Rendering-Pipeline-Prozessoren enthält. Die Rendering-Pipeline blieb zwischen den Cube-Versionen und VolumePro weitgehend konstant. Weiterentwickelt wurden vor allem die Speicherarchitekturen, so konnte der VolumePro 500 erstmals SDRAM-Bausteine verwenden. Die aktuellen VolumePro 1000-Karten arbeiten mit bis zu 4 GByte DDRAM bei 250MHz Takt und maximal zwei vg1000-Bausteinen. Mit dem vg1000 werden noch zusätzlich Oberflächengrafiken unterstützt, die auf der Zusatzgrafikkarte berechnet werden.

Um einen gleichförmigen Datenzugriff zu erreichen, wird das Volumen beim vg500 wie bei Shear-Warp (siehe Kapitel 3.5) in Schichten (Slices) abgetastet, die parallel zu der Würfelseite liegen, die der Projektionsebene am meisten zugewandt ist. Es wird zuerst ein Zwischenbild auf diese Würfelseite abgebildet, das erst in einem zweiten Schritt abgetastet und auf die Projektionsebene entzerrt wird. VolumePro mit dem vg500-Chip übernimmt die Projektion auf die Würfelseite (Shear). Die Entzerrung (Warp) auf die Projektionsebene wird mit Texture-Mapping auf der mitgelieferten Grafikkarte ausgeführt.

Abbildung 4.3 zeigt die Ray-Casting-Pipelines des vg500-Chips, wobei alle Blöcke vier mal vorhanden sind. Aus dem externen Volumenspeicher (Voxel-Memory) wird Schicht um Schicht und Strahl um Strahl jedes Volumelement einmal pro Bild ausgelesen und den Pipelines zugeführt. Da der Volumenspeicher aus vier parallelen Modulen besteht, können die Daten für die vier Pipelines mit dem gleichen Takt von 125MHz ausgelesen und verarbeitet werden.

Ähnlich wie in Kapitel 7.5 beschrieben, wird das Volumen würfelförmig in den SDRAM-Modulen, Bänken und Zeilen aufgeteilt, so dass keine Zugriffskonflikte entstehen können. Im Gegensatz dazu wird ein 2^3 -Subwürfel aber nicht auf verschiedene Module aufgeteilt, sondern in einer SDRAM-Zeile abgelegt, um den 2^3 -Subwürfel im Burst-Modus auslesen zu können. Eine spezielle Schaltung auf dem vg500-Chip sortiert die Volumelemente der 2^3 -Subwürfel je nach Beobachterraichtung in die Schichten und verteilt sie auf die vier Ray-Casting-Pipelines.

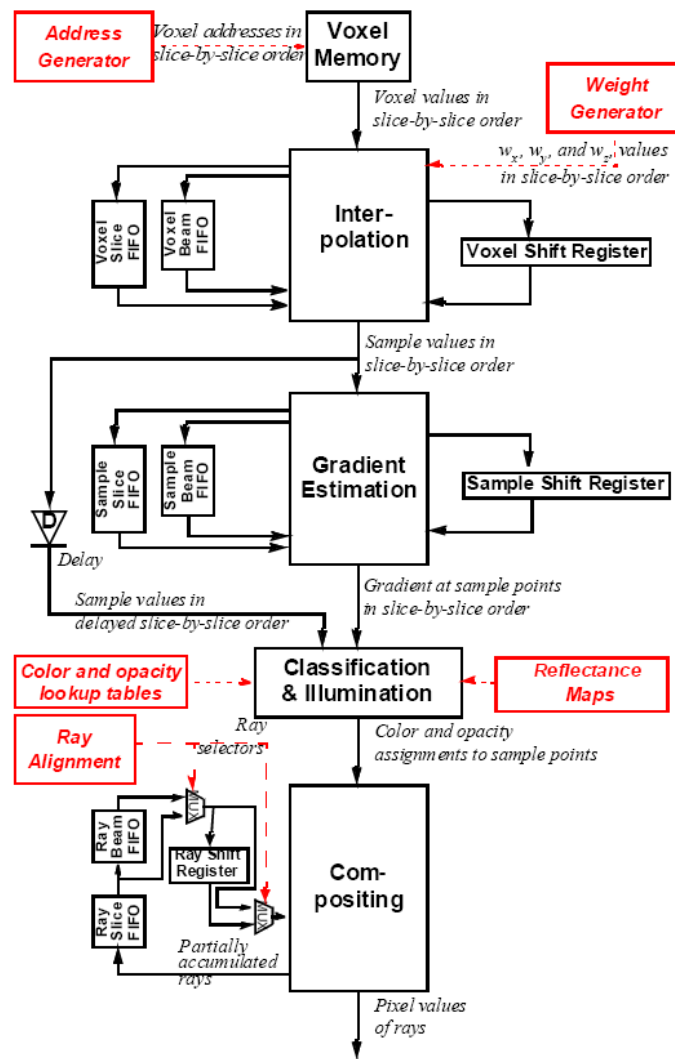


Abbildung 4.3: Ray-Casting-Pipeline von VolumePro (aus [10])

Da alle Volumenelemente nur einmal aus dem Volumenspeicher gelesen werden, müssen die umliegenden Werte für die Interpolation und Gradientenberechnung zwischengespeichert werden. Dies geschieht mit Hilfe der Shift-Register, Slice-Fifos und Beam-Fifos in Abbildung 4.3. Über die Shift-Register findet die Übergabe von Daten an die benachbarte Pipeline statt, wobei die vierte Pipeline wieder mit der ersten verbunden ist.

Damit die Chip-internen Fifos nicht zu groß werden, werden die Schichten in X-Richtung auf 32 Voxel begrenzt (Sectioning). Am Ende eines 32-Voxel Abschnittes werden deshalb die Werte nicht mehr von der vierten Pipeline auf die erste übergeben, sondern in ein Section-Memory zur Zwischenspeicherung geschrieben. Beim Beginn eines neuen Abschnittes liest die erste Pipeline aus dem Section-Memory, anstatt aus dem Shift-Register.

Beim VolumePro1000 wird nicht mehr das Shear-Warp-Verfahren angewandt, sondern eine

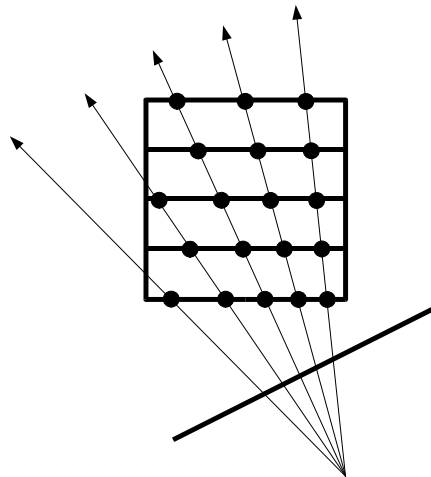


Abbildung 4.4: Shear-Image-Order-Ray-Casting

abgewandelte Form, die sich Shear-Image-Order-Ray-Casting[74] nennt. Mit diesem Verfahren kann die 2D-Interpolation des Warp weggelassen, welche eine zusätzliche Verschlechterung der Bilder nach sich zieht. Erreicht wird dies, indem die Strahlen zwar wie beim Ray-Casting-Verfahren von einem Projektionspunkt durch die Bildpixel laufen, aber die Abtastpunkte nicht äquidistant auf den Strahlen liegen, sondern genau auf den Schichten (Slices) wie bei Shear-Warp (Abbildung 4.4). Dadurch kann die gleiche reguläre Auslesereihenfolge, wie beim VolumePro 500 beibehalten werden. Zusätzlich sind jetzt auch perspektivische Darstellungen und Early-Ray-Termination möglich. Für Space-Leaping wird zwar geworben, es handelt sich allerdings nur um ein Aussparen grober Regionen durch Schnittebenen parallel zu den Seiten des Datenwürfels.

Parallel zur Entwicklung des VolumePro wurden Untersuchungen durchgeführt, die Cube-Architektur für Volume-Ray-Tracing anzupassen [27,28]. Es wird das gleiche Abarbeitungsprinzip, wie beim Race II-System (Kap. 4.6) vorgeschlagen: eine Mischung zwischen Image-Order und Object-Order, wobei versucht wird einmal geladene Blöcke komplett abzuarbeiten. Der Fokus der Untersuchungen liegt in der Abarbeitungsreihenfolge der Speicherblöcke und Cache-Strategien für das Volume-Ray-Tracing-Verfahren. Es wurde ein theoretisches System mit einem oder mehreren Blockprozessoren (abgeleitet vom vg500-Pipeline-Prozessor), einem DSP^I und RDRAM^{II} zugrunde gelegt.

^I DSP: Digitaler Signal Prozessor

^{II} RDRAM: Rambus Dynamic Random Access Memory; pro Kanal 800MB/s Datendurchsatz, 4 Kanäle möglich, siehe auch Kapitel 7.3.1

4.5 Vizard II

VIZARD II ist ein Ray-Casting-System, das in einer Zusammenarbeit der Universität Tübingen und der Philips Research Laboratories in Hamburg entwickelt wurde. Es handelt sich um ein PCI-Board Design, das 256^3 Datensätze für ein 256^2 Bild mit Wiederholraten von 3 bis 7 Bilder pro Sekunde visualisieren kann [44, 45, 46].

Die Hauptbestandteile des Systems sind neben vier getrennten SDRAM-Speicherbänken mit Standard-DIMM-Sockeln, ein DSP zur Startpunkt und Offsetberechnung der Strahlen und ein Xilinx Virtex FPGA, der die Ray-Processing-Unit aus Abbildung 4.5 enthält.

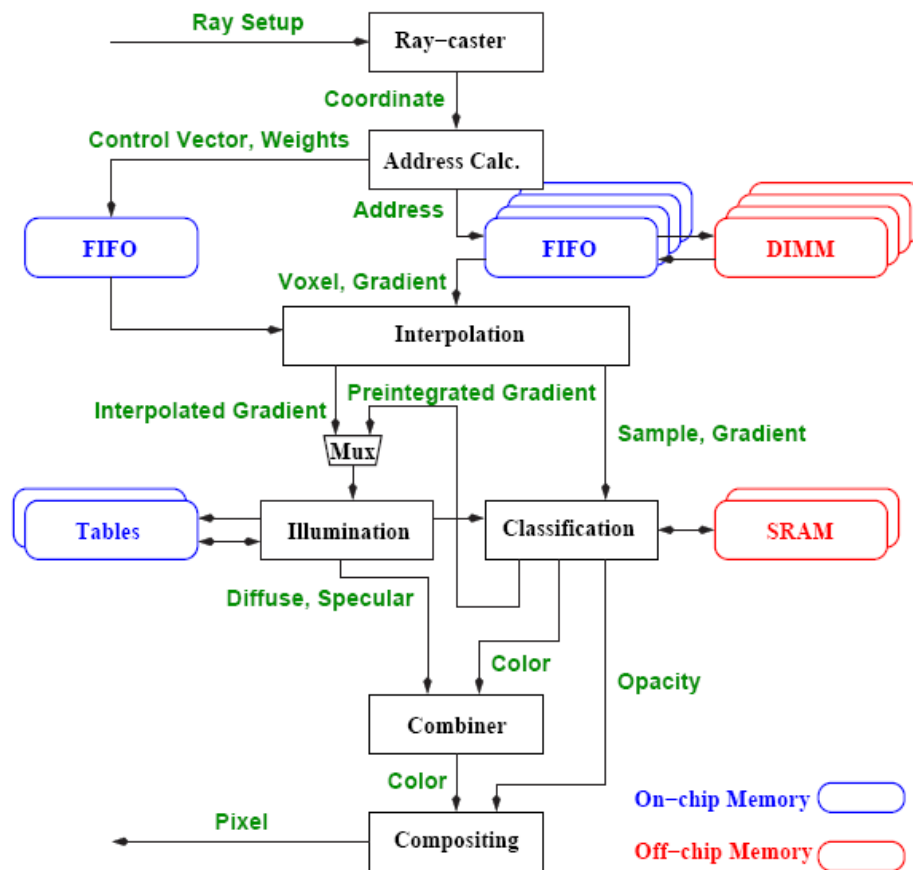


Abbildung 4.5: Aufbau der Ray Processing Unit (RPU) von VIZARD II

Es wurde prinzipiell eine Rendering-Pipeline mit Interpolation, Gradientenberechnung und Shading mit Reflectance-Map wie in Kapitel 6.5 bis 6.7 beschrieben, implementiert. Die von uns in [37] beschriebene Ray-Queue, für einen Multithreading-Ansatz (Kapitel 7.4) wurde mit 9 parallelen Strahlen, allerdings ohne Sortierung, realisiert. Sie wird nur genutzt, um zusätzliche Takte der Adressgenerierung auszugleichen. Early-Ray-Termination findet nur statt, wenn alle 9 Strahlen

unter den eingestellten Schwellwert fallen, weshalb auf die Sortierung wegen der Datenkohärenz verzichtet werden kann. Space-Leaping wurde nicht implementiert.

VIZARD II kommt mit 4 Speicherbänken aus, da in den 64-Bit breiten SDRAM-Modulen die 32 Bit-Volumendaten in Z-Richtung doppelt eingespeichert werden, so dass mit einem Zugriff immer die 8 Volumenelemente eines 2^3 -Subwürfels zur Interpolation ausgelesen werden können.

Materialeigenschaften können über den Dichtewerte adressiert, aus einer Tabelle, in Form von Koeffizienten für ambiente, diffuse und spiegelnde Reflexion, mit berücksichtigt werden.

Einfache Gradienten können während der Visualisierung aus den 2^3 -Subwürfeln errechnet werden. Genauere Gradienten können in einer Vorberechnung bestimmt werden. Das Ergebnis der Vorberechnung wird als Index auf die Einheitsgradienten einer Kugeloberfläche mit den Volumendaten abgespeichert^I. Beim Zugriff auf einen 2^3 -Subwürfel werden somit 8 Verweise auf die 512-Einträge großen Look-Up-Tabellen^{II} mit den Gradienten ausgelesen. Der Gradient am Abtastpunkt wird durch getrennte Interpolation der X,Y,Z-Komponente der 8 Gradienten erreicht.

Zur Vermeidung von Artefakten ist in VIZARD II das Pre-Integrations-Verfahren implementiert. Mit ihm werden Integrationsfehler entlang des Strahls reduziert (siehe Kapitel 5.4).

4.6 Race II

Mit dem Race II-System[12] wurde eine Architektur vorgeschlagen, die versucht die Vorteile von *Object-Order* und *Image-Order* (vgl. Kap 2.1) zu verbinden.

Hierzu wird der Volumenspeicher in 8^3 -Würfel (VAB: Voxel Access Block) unterteilt. Diese werden von vorne nach hinten (*Object-Order*) in einen Zwischenspeicher geladen, von dem aus die weitere Berechnung startet.

Zuerst werden die acht Eckpunkte des VAB auf die Bildebene projiziert (Forward-Projection). Für alle Pixel innerhalb des daraus entstehenden 2D-Abdrucks wird dann ein Image-Order-Ray-Casting-Verfahren innerhalb des VAB durchgeführt (Backward-Projection) (siehe Abbildung 4.6 Step 1 bis 6).

Durch die Zwischenspeicherung der VAB kann der erforderliche Durchsatz des Volumenspeichers gegenüber anderen Image-Order-Systemen stark reduziert werden, da die Daten nur einmal gelesen werden und weitere notwendige Zugriffe nur auf den Zwischenspeicher erfolgen. Dadurch sind auch

I Analog zur Adressierung der Reflexionstabelle in Kapitel 6.7.

II Look-Up-Tabellen (LUT) werden in Hardware durch statische RAM- oder ROM-Speicher realisiert und können komplexe Rechenwerke, beispielsweise für eine mathematische Funktion, ersetzen. Der Parameter der Funktion wird als Adresse an den Speicher gelegt. An der ausgewählten Speicherstelle wird das vorberechnete Ergebnis der Funktion über den Datenbus des Speichers ausgelesen.

größere Umgebungen zur Gradientenberechnung realisierbar.

Durch Doppelpufferung der VAB kann das Laden des nächsten VAB während der Berechnung im Hintergrund stattfinden. Die Berechnungsergebnisse an den Rändern der VAB müssen allerdings zwischengespeichert werden, da sie für die Interpolation und Gradientenberechnung benötigt werden.

Allerdings kann weder Early-Ray-Termination noch Space-Leaping effizient implementiert werden. Da der VAB immer geladen wird, nützt Early-Ray-Termination nur etwas bei Überabtastung. Entsprechend können durch Space-Leaping nur komplette VAB ausgelassen werden, die beim ersten Laden als leer erkannt und markiert werden.

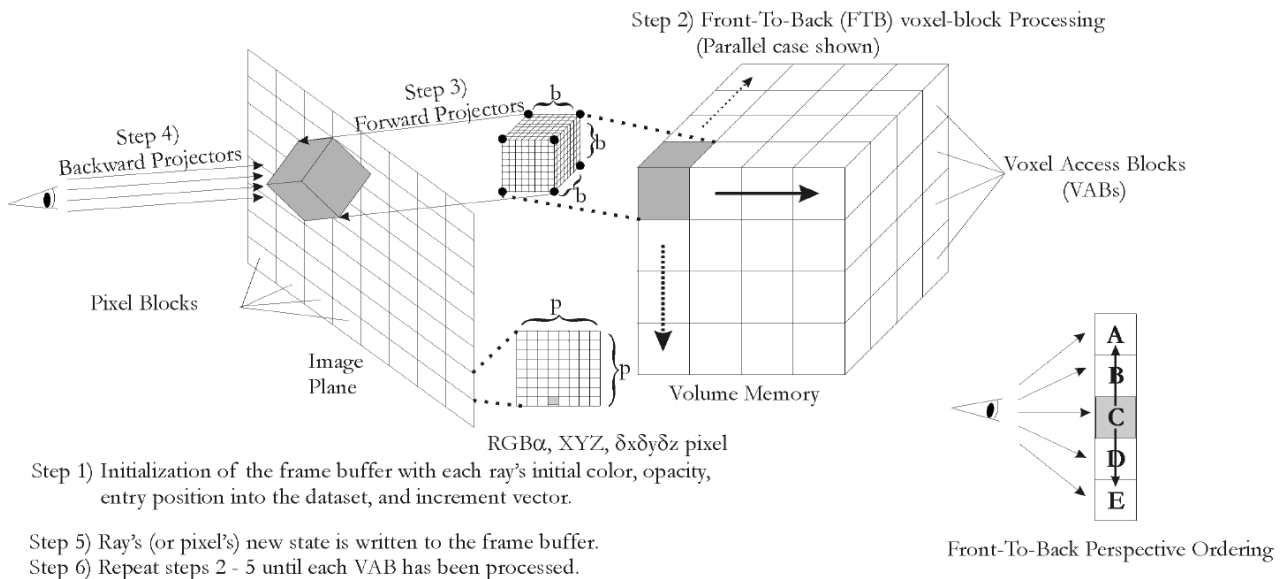


Abbildung 4.6 Algorithmus des Race II-Systems

5 Ray-Casting-Algorithmus

5.1 Übersicht

Im Folgenden soll die Rendering Pipeline, wie sie Levoy [22] beschrieben hat, erläutert werden. In der Reihenfolge der Abarbeitung sind einige Modifikationen entwickelt worden, deren Vor- und Nachteile in den nachfolgenden Abschnitten aufgezeigt werden sollen. Levoy trennt strikt die Abarbeitung der Opazität und der Farbe. Opazität ist ein Maß für die Undurchsichtigkeit – das

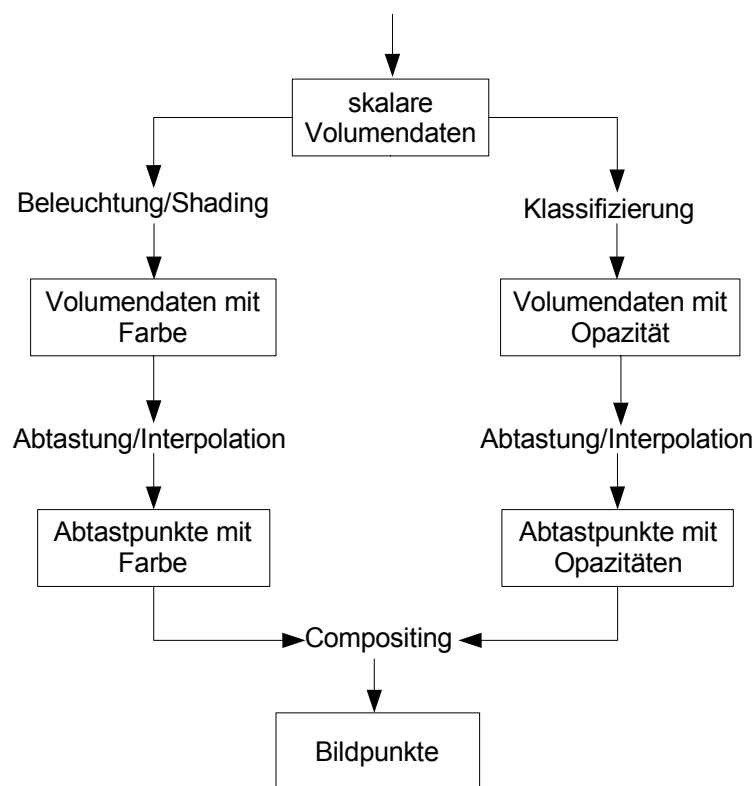


Abbildung 5.1: Rendering Pipeline nach Levoy

Gegenteil zur Transparenz. Sie wird vom Benutzer, bei der Klassifizierung, beliebigen Wertebereichen der Volumendaten zugeordnet. Zum Beispiel können in einer CT-Aufnahme alle Werte unter 1000 auf Null – also völlig transparent – gesetzt werden und darüber linear auf 1 – völlig undurchsichtig – ansteigen. Dadurch sind bei der Visualisierung nur noch die Knochen zu sehen. Die Farbe dagegen wird durch einen Beleuchtungsalgorithmus in Abhängigkeit von lokalen Gradientenvektoren, die als Oberflächennormale interpretiert werden, berechnet. Opazität und Farbe sind somit zusätzliche Werte, die auf zwei getrennten Wegen aus den skalaren Volumendaten gewonnen werden (Abbildung 5.1). Anschließend beginnt bei Levoy der eigentliche

Visualisierungsprozess; es werden Sehstrahlen von jedem Bildpunkt aus durch das Volumen gelegt, und entlang der Strahlen in regelmäßigen Abständen abgetastet. Da die Abtastpunkte normalerweise zwischen den Gitterpunkten der Volumendaten liegen, werden Farbe und die Opazitäten in getrennten Stufen trilinear interpoliert und erst im Compositing miteinander verrechnet, aufaddiert und Farbe und Helligkeit des Bildpunktes bestimmt.

Die in dieser Arbeit verwendete Rendering-Pipeline wurde an einigen Stellen modifiziert. Es wurde zum einen eine zusätzliche Möglichkeit der Klassifizierung nach der Interpolation eingefügt, weshalb man in Bezug auf die Interpolation von Pre- bzw. Post-Klassifikation spricht. Zum Anderen wird der Gradientenvektor aus den Opazitäten gewonnen, also die strikte Trennung zwischen Klassifizierung und Beleuchtung aufgegeben (Abbildung 5.2). Dadurch werden Artefakte vermieden, die entstehen, weil die Gradienten aus den unveränderten Volumendaten gewonnen wurden. Durch die Klassifizierung können die Opazitäten aber derart verändert werden, dass aus

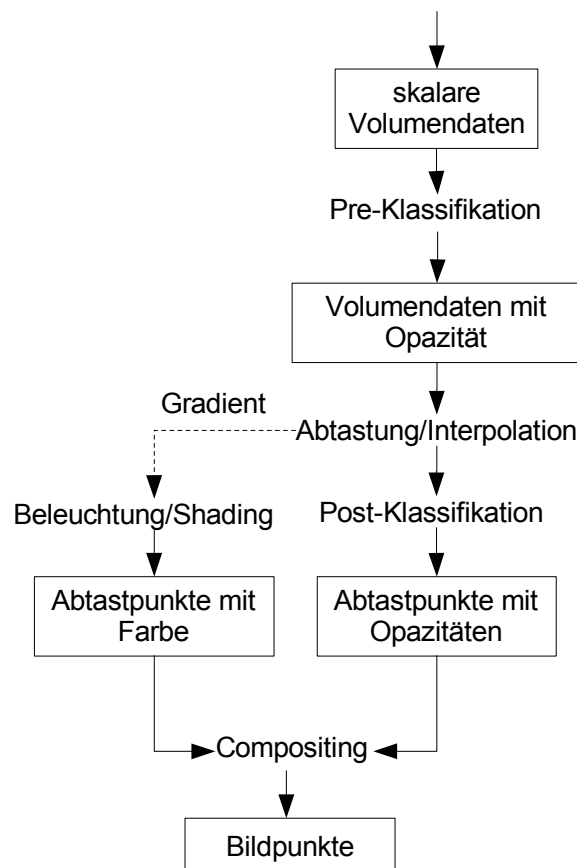


Abbildung 5.2: Modifizierte Rendering

Originaldaten gewonnene Gradienten in die falsche Richtung zeigen. Zum Beispiel können durch Ausblenden bestimmter Dichtewerte neue Kanten entstehen oder verschwinden, die durch die

Gradienten aus den Originaldaten nicht angezeigt werden.

Die Gradienten fallen bei der im Kapitel 6.5 erläuterten Interpolationsschaltung quasi kostenlos ab.

Die einzelnen Stufen der Rendering-Pipeline werden in den folgenden Abschnitten erläutert.

5.2 Klassifizierung

Das Volume-Rendering-Verfahren und damit auch der Ray-Casting-Algorithmus werden zur Visualisierung von dreidimensionalen Skalarfeldern verwendet. Bei der Klassifizierung werden den möglichen Skalarwerten Materialeigenschaften zur Visualisierung zugeordnet. Materialeigenschaften sind beispielsweise die Farbe oder die Opazität. Weitere Möglichkeiten sind die Koeffizienten für das Beleuchtungsmodell, welche die Anteile von ambienter, diffuser und spiegelnder Reflexion festlegen. Sie können durch die Klassifikation jedem Skalarwert unterschiedlich zugewiesen werden.

Möglich ist auch eine Segmentierung der Volumendaten. Hierbei werden unabhängig vom Skalarwert die Volumenelemente einer Gruppierung zugeordnet. Die Segmentierung wird vor dem Visualisieren durchgeführt und beim Render-Vorgang die Segmentnummer ausgewertet. Zum Beispiel die Zugehörigkeit des Volumenelementes zu einzelnen Geweben oder Organen, wie Haut, Knochen oder Muskeln, die der Arzt während der Visualisierung sichtbar oder unsichtbar schalten kann. Dann können die Materialeigenschaften den Gruppierungen und nicht direkt den Skalarwerten zugeordnet sein.

In der hier verwendeten Rendering-Pipeline beschränkt sich die Klassifizierung auf die Zuordnung von Opazitäten und Segmentierung. Farbe und Beleuchtungsparameter werden im Kapitel 8 als Erweiterung behandelt.

Je nachdem, ob die Klassifizierung vor oder hinter der Interpolation durchgeführt wird spricht man von Pre- oder Post-Klassifikation. Die jeweiligen Vor- und Nachteile sollen im Folgenden aufgezeigt werden.

5.2.1 Pre-Klassifikation

Die Pre-Klassifikation ist vor allem zur Darstellung segmentierter Daten notwendig. Die Klassifikation ordnet den Dichtewerten der Originaldaten eine Opazität im Wertebereich zwischen 0 und 1 zu. Der Benutzer kann normalerweise über die Anwendung in einer Diagrammkurve die Funktion der Zuordnung beliebig einstellen.

Da mit diesem Hilfsmittel ganze Objektteile unabhängig von den Dichtewerten ausgeschaltet

werden können, muss dies vor der Interpolation und der Gradientenberechnung geschehen, um die neu entstehenden Kanten richtig darzustellen. Durch die Interpolation kommt es auch zu weichen Übergängen zwischen den Segmenten. Die Interpolation selbst erhält keine Segmentinformation, sie wird durch die Pre-Klassifikation bereits ausgefiltert.

Bei Überabtastung, das heißt, beim Vergrößern des Objektes in Bereiche, in denen das Gitteraster der Originaldaten sichtbar wird, kann im Gegensatz zur Post-Klassifikation das Raster an der Oberfläche würfelförmig sichtbar werden [52].

5.2.2 Post-Klassifikation

Die Post-Klassifikation hat den Vorteil selbst bei starken Vergrößerungen eine glatte Oberfläche zu zeigen, ohne dass das Gitteraster sichtbar wird [52].

Ein Nachteil ist, dass Gradienten die innerhalb der Interpolation berechnet wurden, nicht zu den Opazitäten passen, die durch die Post-Klassifikation entstehen. Es kommt dadurch zu Artefakten. Um diese zu vermeiden, müsste die Gradientenberechnung nachgeschaltet werden. Das macht allerdings sehr große Pufferspeicher notwendig, da die Berechnung alle umliegenden Opazitäten der Nachbarstrahlen benötigt.

5.3 Beleuchtungsmodelle, Shading

Beim Ray-Casting-Algorithmus wird für jeden abgetasteten Punkt entlang eines Sehstrahls eine Reflexion zum Beobachter berechnet. Eine naturnahe Berechnung der Lichtverhältnisse einer Szene kann beliebig komplex werden, wenn beispielsweise Mehrfachreflexionen oder unterschiedliche Materialien berücksichtigt werden. Deshalb verwendet man vereinfachte Modelle, die bei minimaler Rechenzeit ein attraktives Ergebnis liefern. Ein weit verbreitetes Beleuchtungsmodell, das diese Eigenschaften besitzt, wurde bereits 1975 von Phong für Oberflächengrafik entwickelt, ist in [32] beschrieben und wird oft als Phong-Shading bezeichnet.

Im Folgenden werden die Bestandteile dieses Beleuchtungsmodells beschrieben.

Bei Farbbildern sind die Berechnungen für jeden Farbkanal Rot, Grün, Blau getrennt zu berechnen. Die Koeffizienten, die den Beitrag der einzelnen Reflexionsanteile bestimmen, können pro Farbkanal unterschiedlich sein.

5.3.1 Ambiente Reflexion

Die ambiente Reflexion ist der einfachste Bestandteil des Phong-Shading; er ist unabhängig von der Richtung, in welcher die Lichtquelle oder der Beobachter sich befinden.

Die ambiente Reflexion I_a wird wie folgt berechnet:

$$I_a = I_0 k_a$$

wobei I_0 die Intensität der Lichtquelle an der Oberfläche und k_a als Koeffizient den Anteil der Reflexion angibt und im Wertebereich zwischen 0 und 1 liegt.

Man benötigt diesen Bestandteil der Reflexion, um Objektteile, die ansonsten durch Abschattung nicht beleuchtet sind, noch zu erkennen. Es gibt jedoch keinen direkten physikalischen Bezug für diesen Lichtbestandteil. Er kann aber als Approximation von Mehrfachreflexionen gesehen werden. Ambiente Reflexion alleine leuchtet die Objekte völlig gleichmäßig aus, so dass keine Struktur erkennbar ist.

5.3.2 Diffuse Reflexion

Die diffuse Reflexion beleuchtet das Objekt heller, je mehr die Oberfläche in die Richtung der Lichtquelle zeigt. Sie ist unabhängig von der Position des Beobachters und simuliert das Verhalten von matten, nicht spiegelnden Oberflächen. Die Intensität ist proportional zum Cosinus des Winkels zwischen dem Vektor zur Beleuchtungsquelle \vec{L} und der Oberflächennormalen \vec{N} . Über den Koeffizienten k_d wird der Anteil der diffusen Reflexion eingestellt. I_0 gibt wieder die Intensität der Lichtquelle an:

$$I_d = I_0 k_d \cos \theta \quad \text{mit} \quad \cos \theta = \vec{N} \cdot \vec{L} \quad \text{ergibt sich}$$

$$I_d = I_0 k_d (\vec{N} \cdot \vec{L})$$

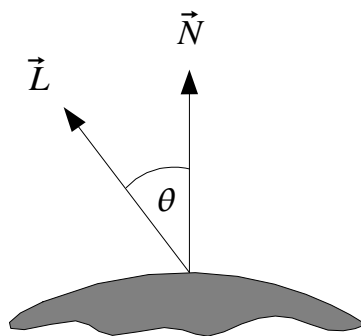


Abbildung 5.3: Diffuse Reflexion

5.3.3 Spiegelnde Reflexion

Helle Lichtspots auf glatten, spiegelnden Flächen, wie Plastik oder Metalloberflächen werden durch diesen Anteil simuliert. Die Spiegelnde Reflexion ist einmal vom Winkel θ zwischen der

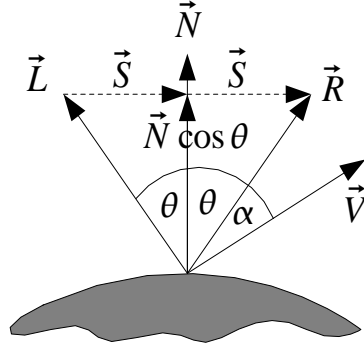


Abbildung 5.4: Spiegelnde Reflexion

Beleuchtungsquelle \vec{L} und der Oberflächennormalen \vec{N} abhängig, aber auch vom Winkel α zwischen dem Beobachtungsvektor \vec{V} und dem Reflexionsvektor \vec{R} , der sich geometrisch durch die Rechnung "Einfallswinkel gleich Ausfallswinkel" des Lichtes an der Oberfläche berechnet. Die Größe des Lichtspots wird durch die Potenz n des Cosinus von α bestimmt. Für eine ideale Spiegeloberfläche geht n gegen unendlich. Auch der Betrag dieses Anteils wird durch einen Koeffizienten k_s und der Intensität der Lichtquelle I_0 gesteuert.

$$I_s = I_0 k_s \cos^n \alpha$$

$$I_s = I_0 k_s (\vec{R} \cdot \vec{V})^n$$

Berechnung von \vec{R} nach Abbildung 5.4:

$$\vec{R} = \vec{N} \cos \theta + \vec{S}$$

$$\vec{S} = \vec{N} \cos \theta - \vec{L} \text{ daraus folgt}$$

$$\vec{R} = 2 \vec{N} \cos \theta - \vec{L} \text{ mit } \cos \theta = \vec{N} \cdot \vec{L} \text{ lässt sich } \vec{R} \text{ berechnen als}$$

$$\vec{R} = 2 \vec{N} (\vec{N} \cdot \vec{L}) - \vec{L}$$

Für den spiegelnden Lichtanteil mit $n=1$ ergibt sich:

$$I_s = I_0 k_s (2 \vec{N} (\vec{N} \cdot \vec{L}) - \vec{L}) \cdot \vec{V}$$

5.4 Abtastung, Interpolation und Integration

Medizinische Aufnahmen durch CT oder MRT haben nur eine Auflösung von etwa einem Millimeter. Kleinere Strukturen sind somit aus den Daten nicht rekonstruierbar. In Abbildung 5.5 sieht man als Beispiel deutlich, dass die feine Spitze im Maximum der Kurve in den diskreten Volumendaten nicht mehr zu sehen ist.

Ähnliches gilt auch für künstlich erzeugte Volumendaten für die Computergrafik, die immer nur eine endliche Ortsauflösung besitzen.

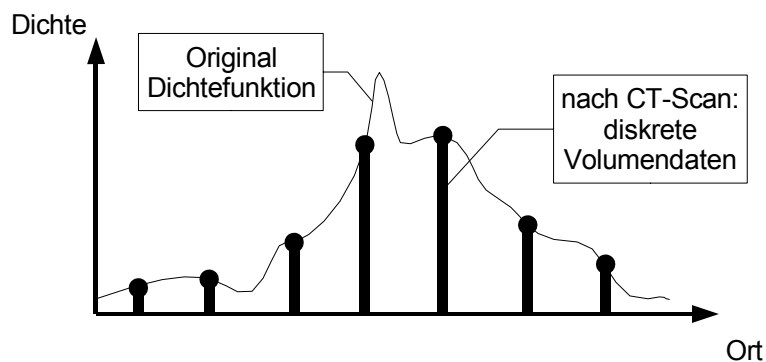


Abbildung 5.5: Erzeugung der Volumendaten durch einen CT-Scan

Als Eingangsdaten für den jeweiligen Volume-Rendering-Algorithmus stehen nur die diskreten Volumendaten auf einem orthogonalen Gitternetz, wie in Abbildung 5.6 angedeutet zur Verfügung. Dadurch kann aus den Volumendaten auch keine detailliertere Information gewonnen werden.

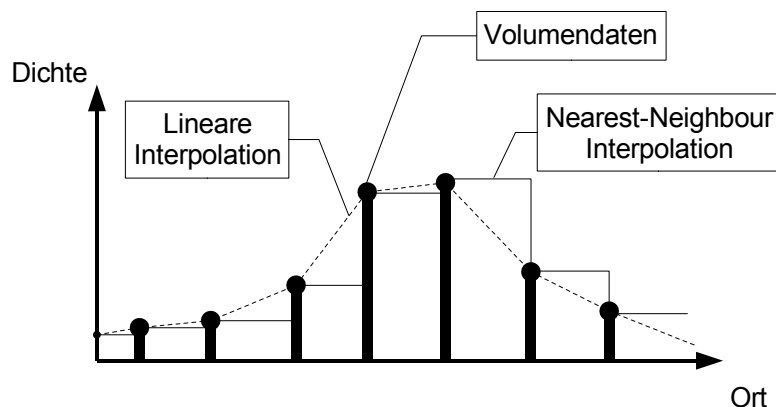


Abbildung 5.6: Abtastung der diskreten Volumendaten und Interpolation

Bei der Visualisierung mit dem Ray-Casting-Algorithmus werden von einer beliebig im Raum liegenden Projektions- oder Bildebene, virtuelle Sehstrahlen durch das Volumen gesendet. Diese

können perspektivisch auseinander laufen. Entlang der Strahlen werden in gleichmäßigen Abständen Abtastpunkte gesetzt und deren Werte an die weitere Berechnung weitergegeben. Da die Abtastpunkte zwischen den Gitterpunkten der Volumendaten liegen können, muss man sie entweder einem Gitterpunkt zuordnen, oder durch Interpolation einen Mittelwert berechnen.

Die Zuordnung zu einem Gitterpunkt, die auch Nearest-Neighbour-Interpolation genannt wird, entstehen Treppenstufen (Abbildung 5.6). Sie stellt eine eher schlechte Approximation der originalen Dichtefunktion dar. Sie hat den Vorteil, dass nur ein Wert ausgelesen werden muss, führt aber zu einer sehr schlechten Bildqualität.

Für den aus der Signaltheorie bekannten Idealfall, die Interpolation mit Hilfe einer $\frac{\sin x}{x}$ -Approximationsfunktion oder dem nahe kommenden Filter, müsste eine größere Nachbarschaft zur Berechnung ausgelesen werden, vor allem da die Interpolation in drei Dimensionen durchgeführt werden muss. Da bei Hardware-Realisierungen für Volume-Rendering das Auslesen der Daten den Flaschenhals bilden, kommt dies aber für Echtzeit-Systeme nicht in Frage. Der beste Kompromiss zwischen dem Aufwand für das Auslesen der Umgebung und der Bildqualität ist sicherlich die lineare, oder im dreidimensionalen, eine trilineare Interpolation. Sie benötigt nur die acht um den Abtastpunkt liegenden Datenpunkte.

Die trilineare Interpolation zeigt Artefakte, die vor allem bei starken Vergrößerungen als Kanten an den Rändern der Volumenelemente sichtbar werden. Es gibt Ansätze, die Oberfläche beispielsweise durch eine kubische Interpolation zu glätten [53], allerdings bedeutet das einen wesentlich größeren Rechenaufwand pro Abtastpunkt und erfordert ebenfalls das Einlesen einer größeren Umgebung. Glatte Oberflächen sind für die Computergrafik interessant, in der medizinischen Anwendung ist es aber besser, wenn die Grenzen der Auflösung bei der Visualisierung noch erkennbar bleiben.

Wenn die trilineare Interpolation durch die Hardware vorgegeben ist, was geschieht dann bei Überabtastung der diskreten Volumendaten?

Wie Abbildung 5.7 verdeutlichen soll, bekommt man durch die Interpolation zwar zusätzliche Zwischenwerte, aber es wird keine zusätzliche Information gewonnen. Will man also alle in den Volumendaten enthaltene Details erkennen, würde eine Abtastweite, die dem Gitterraster entspricht, ausreichen. Vor allem reduziert eine Überabtastung erheblich die Bildwiederholrate, da eine Verdopplung der Abtastrate im Raum acht mal so viele Abtastpunkte generiert.

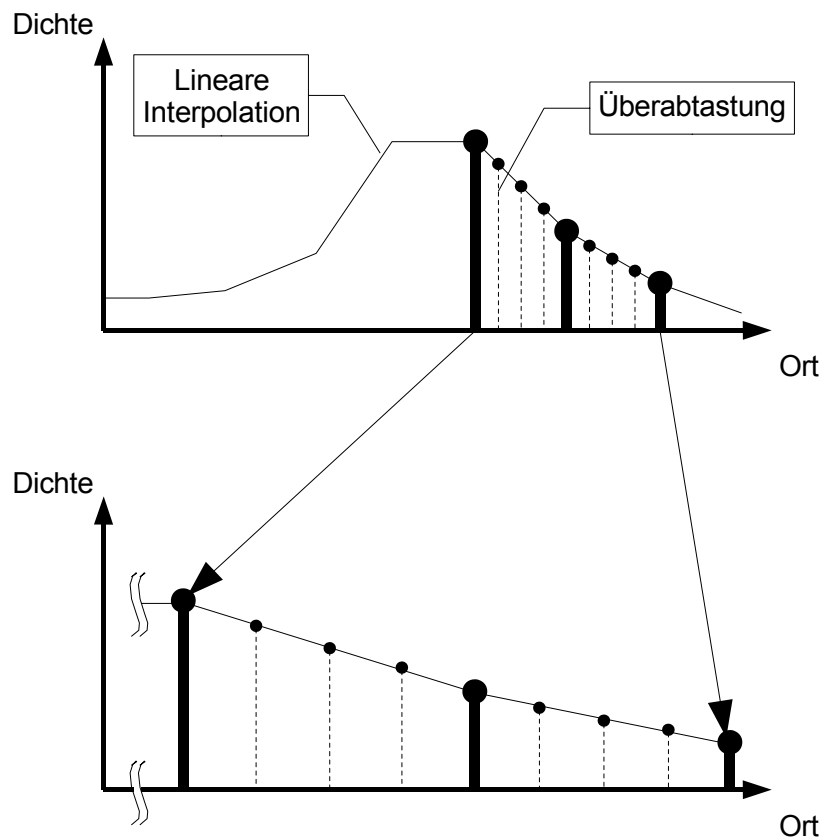


Abbildung 5.7: Keinen Informationsgewinn durch Überabtastung bei linearer Interpolation

Es gibt dennoch einige Gründe, die für die Überabtastung sprechen. Einmal wird sie oft angewendet, um eine Vergrößerung eines Ausschnitts zu erhalten, oder um bei perspektivischer Darstellung auch im hinteren Volumenbereich noch alle Volumenelemente mit einzubeziehen.

Ein weiterer Grund, der eine Überabtastung sinnvoll macht, ist eine Verringerung des Fehlers bei der numerischen Integration der Reflexionen entlang des Strahls. In den meisten Realisierungen wird der Einfachheit halber die numerische Integration mit der Rechteckregel verwendet, also alle Reflexionen werden einfach addiert. Wie *de Boer* in [82] beschreibt, ist der dabei entstehende Fehler nicht unerheblich. Er zeigt sich als Artefakt mit gut sichtbaren regelmäßigen Streifen und kann durch Überabtastung minimiert werden.

Ein anderer Weg sind verbesserte Integrationsverfahren, wie zum Beispiel die lineare Approximation mit der Trapezregel durch Pre-Integration. Dadurch können auch ohne Überabtastung gute Bilder gewonnen werden. Es erfordert aber mehr Berechnungsschritte pro Abtastpunkt und macht als Hardware-Realisierung das Rechenwerk aufwändiger [47]. Das Verfahren kann ohne weiteres auch in der hier vorgestellten Architektur integriert werden. Es realisiert das Trapezverfahren mit Hilfe von Look-Up-Tabellen. Das Shading basiert hier nicht, wie

bei der herkömmlichen Methode, auf dem Wert eines Abtastpunktes, sondern zweier aufeinander folgenden Abtastpunkten. Mit beiden Abtastpunkten werden innerhalb des Shadings und der Klassifikation die Look-Up-Tabellen adressiert. In [44] wurde beispielsweise eine Gewichtung für die Korrektur der Gradienten, ein Absorptionsfaktor und die Faktoren für den ambienten, spiegelnden und diffusen Lichtanteil, berücksichtigt. In [78] wurde das Verfahren näher untersucht. Der Autor konnte zeigen, dass durch Pre-Integration alleine nicht alle Artefakte verschwinden. Erst eine zusätzliche vierfache Überabtastung konnte einwandfreie Ergebnisse hervorbringen. Eine noch höhere Überabtastung zeigte keine weitere Verbesserung.

Eine weiter interessante Methode die Dichtefunktion mittels Bézier-Kurven nachzubilden wurde in [25] veröffentlicht, so dass ebenfalls ohne Überabtastung artefaktfreie Bilder entstehen.

5.4.1 Absorption

Beim ursprünglichen Phong-Shading-Beleuchtungsmodell, das für Oberflächengrafiken entwickelt wurde, ist für den spiegelnden und den diffusen Anteil eine Dämpfung der Lichtquelle vorgesehen, die abhängig von der Entfernung zur Lichtquelle ist. Dadurch erscheinen Flächen, die näher an der Lichtquelle liegen, heller. Dies ist für Volumengrafiken nicht praktikabel, da hier Flächen, die näher am Beobachter liegen, hervorgehoben werden sollen. Deshalb wird eine Absorption des Lichtes beim Durchgang durch ein Volumenelement definiert, die proportional zur Opazität des Volumenelementes ist. Je größer also die Opazität ist, um so größer ist die Absorption des Lichtes.

Beim Ray-Casting-Verfahren wird die Absorption beim Durchgang durch die Volumenelemente nur für die Reflexionsanteile zum Beobachter berücksichtigt. Die Lichtquellen beleuchten jedes Volumenelement gleich. Die Reflexion nach dem Durchgang durch ein Volumenelement i wird näherungsweise [22, 32] berechnet durch:

$$I_{i-1} = I_i * (1 - O_i)$$

Die Opazität O wurde durch die Klassifikation auf einen Wert zwischen 0 und 1 skaliert. Bei der maximal möglichen Opazität wird das Licht somit schon bei einem Volumenelement vollständig absorbiert.

Zur Modellierung verschiedener Materialeigenschaften kann die Opazität für die einzelnen Farben des Lichtes unterschiedlich sein.

5.5 Compositing

In der Oberflächengrafik wird mit Compositing die Überlagerung mehrerer hintereinander liegender

Bilder oder Objekte zu einem Gesamtbild beschrieben [32]. In der Volumengrafik werden die Abtastpunkte entlang eines Lichtstrahls als hintereinander liegende, halbtransparente Bildpunkte betrachtet, die auf gleiche Weise überlagert werden.

Für jeden Abtastpunkt wird zuerst die Gesamtreflexion durch Addition der ambienten, diffusen und spiegelnden Lichtanteile bestimmt:

$$I_{phong} = I_a + I_d + I_s$$

Für die Reflexion des vorhergehenden Abtastpunktes wird die Absorption berechnet und die Reflexion im aktuellen Abtastpunkt wird mit der Opazität gewichtet. Beide werden addiert und ergeben die Gesamtreflexion in Richtung des Beobachters:

$$I_{out,i} = I_{phong,i} O_i + I_{out,i-1} (1 - O_{i-1})$$

Die Integration über alle Abtastpunkte sieht wie folgt aus:

$$I_{ges} = \sum_{i=0}^N \left[I_{phong,i} O_i \prod_{j=0}^{i-1} (1 - O_j) \right] \text{ wobei } I_{phong,N} = \text{Hintergrund} \text{ und } O_N = 1, \text{ wodurch der}$$

hinterste Abtastpunkt als undurchsichtig angenommen wird.

Es werden über alle N Abtastpunkte eines Sehstrahls die Reflexionen addiert ($\sum_{i=0}^N \dots$), wobei die

Reflexionen, durch die Absorption, beim Durchgang durch das Volumen bis zur Projektionsebene

abgeschwächt werden ($\prod_{j=0}^{i-1} \dots$).

Vor allem bei medizinischen Anwendungen ist man an der Visualisierung von Übergängen zwischen Geweben interessiert und nicht an den homogenen Bereichen innerhalb.

Um die Übergänge mehr zu betonen kann jede Opazität eines Abtastpunktes mit dem Betrag des Gradienten $|\nabla f_i|$ an dieser Stelle gewichtet werden. Der Betrag des Gradienten kann als Wahrscheinlichkeitsmaß für die Existenz einer Oberfläche angesehen werden:

$$I_{ges} = \sum_{i=0}^N \left[I_{phong,i} \cdot O_i \cdot |\nabla f_i| \prod_{j=0}^{i-1} (1 - O_j) \right]$$

Levoy [22] macht diese Gewichtung schon während der Klassifizierung und rechnet den Gradienten schon in die Opazität mit ein. Für die Realisierung in Hardware ist dieser Schritt beim Compositing einfacher, da die Klassifizierung durch eine einfache Look-Up-Tabelle (Dichte zu Opazität) umgesetzt werden kann. Eine Look-Up-Tabelle mit der Dichte und den drei Komponenten des Gradienten als Eingangsvektor wäre schwerer realisierbar.

6 Hardware-Implementierung der Rendering-Pipeline

6.1 Pipeline-Architektur

Für den Ray-Casting-Algorithmus bauen die Rechenschritte entlang eines Strahls aufeinander auf, da das Ergebnis einer Berechnung für die folgende benötigt wird. Man kann deshalb die Berechnung nicht beschleunigen, in dem man mehrere Rechenwerke parallel verwendet. Zur Beschleunigung muss die Berechnung deshalb in einer Pipeline-Architektur erfolgen.

In einer Pipeline [32] wird die Berechnung in einzelne, synchron ablaufende Stufen unterteilt und durch Register, als Zwischenspeicher, getrennt. In jeder Stufe werden in einer kombinatorischen Schaltung Berechnungen durchgeführt und mit einem Taktsignal das Ergebnis in Register übernommen. Diese Register geben das Ergebnis während des nächsten Taktzykluses an die nächste Pipelinestufe weiter. Die maximale Taktfrequenz orientiert sich an der langsamsten Stufe, weshalb die einzelnen Stufen etwa die gleiche Laufzeit haben sollten.

Da alle Stufen gleichzeitig arbeiten, spricht man bei einer Pipeline von einer zeitlichen Parallelisierung des Arbeitsablaufes.

Kenngrößen einer Pipeline sind Durchsatz und Verzögerung. Der Durchsatz entspricht der Gesamtrate, mit der die Daten verarbeitet werden und die Verzögerung ist die Zeit, die ein Datum für den kompletten Durchlauf benötigt. Falls der Durchsatz nicht ausreicht, müssen die Pipeline-Stufen, soweit möglich, weiter unterteilt werden und zusätzliche Verzögerungsstufen eingebaut werden.

Dies soll am Beispiel einer linearen Interpolationsberechnung, wie sie in Kapitel 6.5 benötigt wird, erläutert werden [34].

Die lineare Interpolation (Abbildung 6.1) hat die Form $c = bx + a(1 - x)$. In der Grundform sind eine Addition, eine Subtraktion und zwei Multiplikationen erforderlich. Eine einfache Umstellung

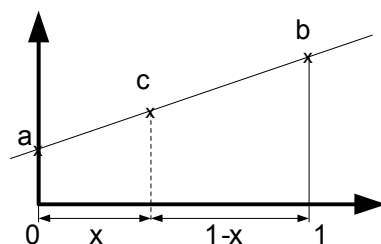


Abbildung 6.1: lineare Interpolation

zu $c = a + (b - a) \cdot x$ spart eine Multiplikation.

Abbildung 6.2 zeigt die zugehörige Schaltung, als eine Pipelinestufe. Angenommen, eine Addition benötigt in der Zielhardware 5 ns und eine Multiplikation 15 ns, dann wäre die Gesamtlaufzeit 25 ns. Diese Stufe könnte mit 40 MHz betrieben werden und hätte nur eine Taktperiode Verzögerung.

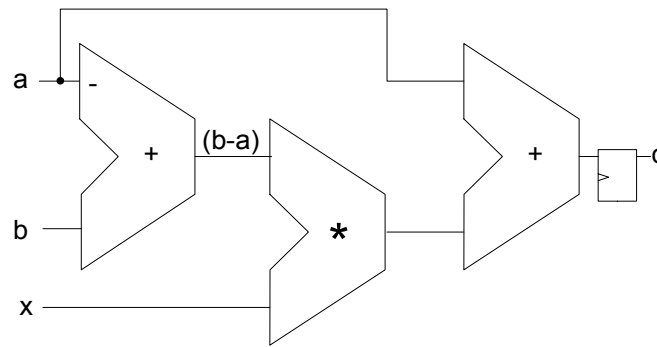


Abbildung 6.2: Schaltung der linearen Interpolation

Soll eine höhere Taktfrequenz erreicht werden, kann man dies durch Einfügen zusätzlicher Register ermöglichen (Abbildung 6.3). Die maximale Taktfrequenz der Schaltung wird durch die langsamste Stufe, in diesem Fall der Multiplizierer mit 15 ns, auf 66 MHz begrenzt. Man nimmt allerdings jetzt in Kauf, dass das zum Eingang passenden Ergebnis erst drei Takte später ausgelesen werden kann; trotzdem wurde der Datendurchsatz entsprechend der Taktfrequenz erhöht.

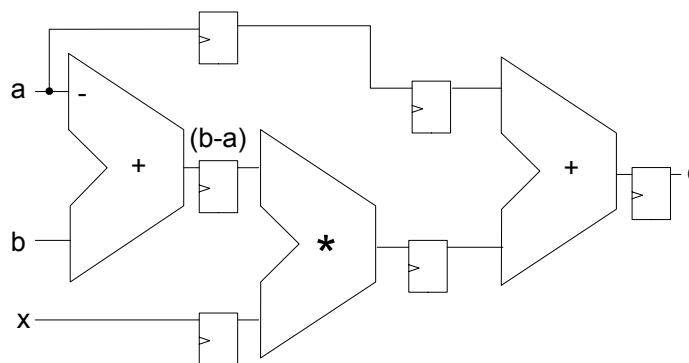


Abbildung 6.3: Schaltung der linearen Interpolation als Pipeline

Alle im Folgenden erläuterten Berechnungsstufen für das Volume-Rendering-Verfahren müssen in dieser Hinsicht auf die erforderliche Taktfrequenz unter Berücksichtigung des zur Verfügung stehenden Platzes optimiert werden. Da dies stark von der Ziel-Hardware und von den dort zur Verfügung stehenden Grundelementen abhängt, werden nur Prinzipschaltbilder ohne Zwischenregister gezeigt.

6.2 Übersicht über die Rendering-Pipeline

In Abbildung 6.4 ist der Ablaufplan der Rendering-Pipeline zu sehen. Die Hardware-Realisierung der einzelnen Komponenten wird im Folgenden beschrieben.

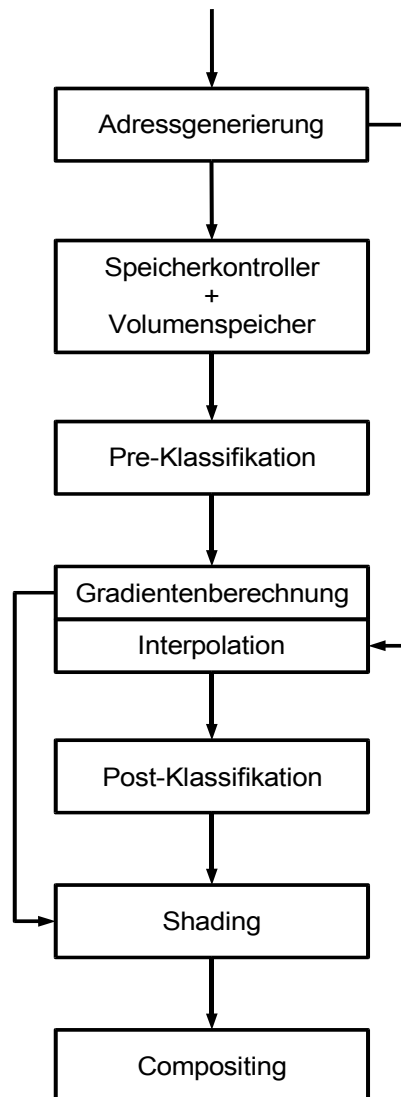


Abbildung 6.4: Verwendete Rendering-Pipeline

Da in dieser Arbeit das Augenmerk auf der Realisierbarkeit der algorithmischen Optimierungen lag und die Implementierungsmöglichkeit in einer FPGA-Hardware im Vordergrund stand, wird eine möglichst einfache Rendering-Pipeline gewählt. So wurde auf die Implementierung von Farbe verzichtet und die Datenbusbreiten auf ein vertretbares Minimum reduziert und durch die Simulation überprüft.

Alle Pipeline-Stufen und der Volumenspeicher arbeiten mit dem gleichen Takt, um einen

maximalen Datendurchsatz zu erhalten. Der Speicherkontroller kann über ein Steuersignal die gesamte Abarbeitung unterbrechen, wenn Seitenwechsel oder Refresh-Zyklen der Speicher notwendig werden.

6.3 Adressgenerierung

Die Adressgenerierung muss die Volumenspeicheradressen der acht Volumenelemente eines 2^3 -Subwürfels (S_0 - S_7) um einen Abtastpunkt bestimmen und an die Speicherkontroller weitergeben. Die Position des Abtastpunktes zwischen den Gitterpunkten der Volumenelemente wird an die Interpolation weitergegeben (siehe Abbildung 6.5). Dies muss fortlaufend entlang eines Sehstrahls für jeden Abtastpunkt ablaufen.

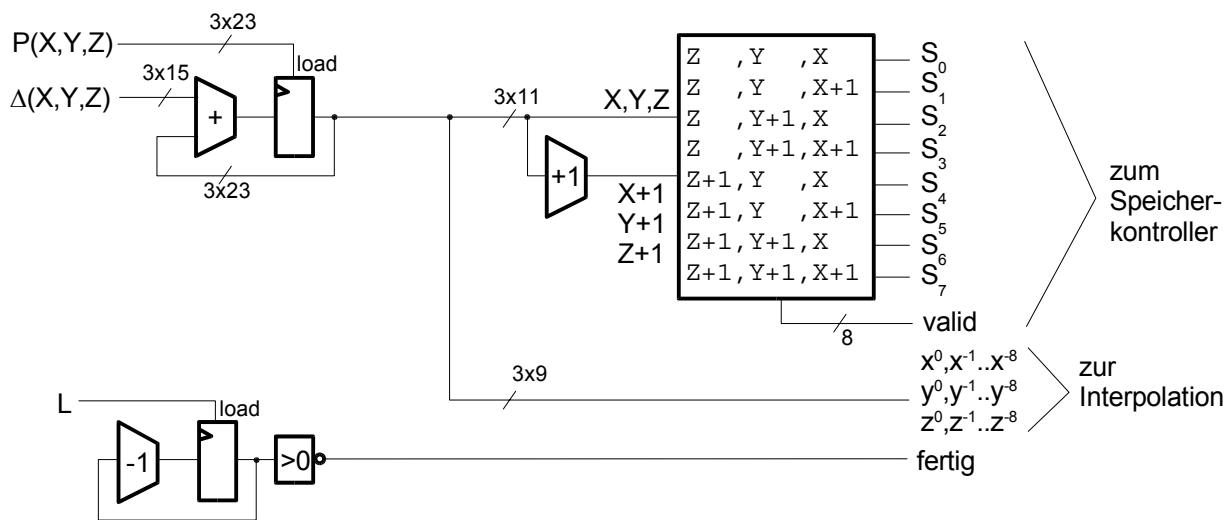


Abbildung 6.5: Prinzipschaltbild des Adressengenerators für einen Strahl

Die Adressgenerierung erhält als Parameter einen dreidimensionalen Startpunkt $P(X,Y,Z)$ des Sehstrahls und einen Differenzvektor $\Delta(X,Y,Z)$, der die Richtung des Sehstrahls und dessen Betrag die Abtastweite angibt. Zusätzlich wird noch die Gesamtlänge des Sehstrahls (L) benötigt.

Die drei Koordinaten des Startpunktes und des Differenzvektors können durch binäre Festkommazahlen mit Vorzeichen repräsentiert werden. So wird beispielsweise $P(X,Y,Z)$ dargestellt mit den Bit:

	<i>Vorzeichen</i>	<i>Vorkomastellen</i>	<i>Nachkomastellen</i>
<i>X-Koordinate</i>	v_x	$x^n..x^0$	$x^{-1}..x^{-m}$
<i>Y-Koordinate</i>	v_y	$y^n..y^0$	$y^{-1}..y^{-m}$
<i>Z-Koordinate</i>	v_z	$z^n..z^0$	$z^{-1}..z^{-m}$

Die aufeinander folgenden Abtastpunkte werden durch wiederholte Addition des Differenzvektors zum Startpunkt bestimmt. Mit jeder Addition wird der Längenparameter des Sehstrahls dekrementiert, um so ein Abbruchkriterium zu erhalten.

Die Vorkommastellen der Koordinaten des Abtastpunktes bestimmen das erste Volumenelement im Volumenspeicher. In Abbildung 6.5 entspricht dies dem Volumenelement S_0 . Die anderen 7 Volumenelemente werden durch entsprechende Erhöhung der X-,Y-und Z-Koordinaten bestimmt. Die räumliche Anordnung der Volumenelemente ist in Abbildung 6.6 zu sehen.

Die einzelnen Bit der Vorkommastellen können direkt den Adressbit des Volumenspeichers, wie in Kapitel 7.5 erläutert, zugewiesen werden.

Die drei Nachkommaanteile ($x^{-1}..x^{-m}$, $y^{-1}..y^{-m}$, $z^{-1}..z^{-m}$) entsprechen der Position des Abtastpunktes zwischen den Gitterpunkten des Volumenrasters und werden, zumindest teilweise (8 Bit), an die Interpolation weitergereicht. Zusätzlich werden noch die LSB^I der Vorkommastellen (x^0 , y^0 , z^0) in der Interpolation benötigt, was in Kapitel 6.5 erläutert wird.

Die Schaltung der Adressgenerierung besteht prinzipiell nur aus Registern, Addierern und einigen Vergleichern.

Die Vergleiche überprüfen, ob der Abtastpunkt sich innerhalb des Volumens befindet und generieren für jede Adresse ein *valid*-Signal. An den Randgebieten und außerhalb des Volumens können mehrere Koordinaten der acht Volumenelemente ungültig sein. Falls die Anzahl der Vorkommastellen genau auf den Speicher abgestimmt sind, reicht die Überprüfung, ob die Koordinate negativ ist, da dies bei einem Übertrag der Addition ebenfalls vorkommt. Ansonsten müsste auch auf die maximale Dimension des Speichers überprüft werden.

Falls eine der Koordinaten eines Volumenelementes ungültig ist, dürfen keine Daten aus dem Speicher gelesen werden, sondern es muss eine 0 statt eines Speicherwertes zur Berechnung in der Interpolation verwendet werden.

Genauigkeitsabschätzung:

Die Busbreiten wurden wie folgt berechnet.

Ein Volumen mit der Größe 1024^3 soll maximal zweifach überabgetastet werden. Dabei soll der maximale Fehler beim Durchgang durch das gesamte Volumen etwa ein halbes Volumenelement betragen.

Bei zweifacher Überabtastung ist die doppelte Anzahl an Additionen zur Durchquerung des Würfels

I LSB: Least Significant Bit, Bit mit der geringsten Wertigkeit

pro Koordinate notwendig. Ohne die Diagonale des Würfels zu berücksichtigen, wären dies 2048 Additionen.

Soll dabei der Fehler maximal ein halbes Volumenelement betragen, muss die kleinste darstellbare

Schrittweite $\frac{0.5}{2048} = 244 \cdot 10^{-6}$ betragen. Daraus kann für die Adressberechnung minimale

Genauigkeit als Anzahl der binären Vor- und Nachkommastellen berechnet werden:

Vorkommastellen: $\log_2(1024) = 10$

Nachkommastellen: $|\log_2(244 \cdot 10^{-6})| = 12$

Die hohe Genauigkeit der Nachkommastellen ist nur für die rekursiv ausgeführte Addition des Differenzvektors notwendig. Für die Weitergabe an die Interpolation reichen 8 Bit aus.

Die Koordinaten des Startpunktes haben somit eine Breite von 23 Bit, inklusive Vorzeichen.

Beim Differenzvektor genügen 2 Vorkommastellen und ein Vorzeichenbit, da selten mehr als eine vierfache Unterabtastung sinnvoll ist, daher werden für seine Koordinaten 15 Bit benötigt.

6.4 Speichercontroller

Der Speichercontroller wird hier nur der Vollständigkeit halber im Vorgriff auf Kapitel 7.3.2 erwähnt, in welchem die Funktionsweise von SDRAM^I-Bausteinen und die daraus resultierende Arbeitsweise des Controllers beschrieben werden.

Für die Realisierung eines Speichercontrollers für SDRAM-Bausteine in FPGAs gibt es für die größeren FPGA-Bausteine fertige Makrozellen. Diese unterstützen aber nur ein Speichermodul. Für das hier beschriebene Ray-Casting-System sind aber acht getrennt adressierbare Speichermodule notwendig, um die acht Volumenelemente parallel für die Interpolation zur Verfügung stellen zu können.

Es wäre möglich, für die Speichercontroller einfach acht Makrozellen zu verwenden. Ressourcen schonender ist es allerdings, einen speziellen Controller für acht Speichermodule zu entwickeln, da viele Funktionen, wie der Refresh, auf alle Module gleichzeitig angewendet werden können. Nur die Funktionen für Lesen und Schreiben müssen für die acht Module getrennt ausgelegt werden. Die genaue Realisierung des Speichercontrollers, als VHDL-Beschreibung, ist in [33] einzusehen.

6.5 Interpolation und Gradientenbestimmung

Wie in Kapitel 5.4 erläutert, ist der beste Kompromiss zwischen Aufwand und Bildqualität bei der

I SDRAM: Synchronous Dynamic Random Access Memory

Abtastung die trilineare Interpolation.

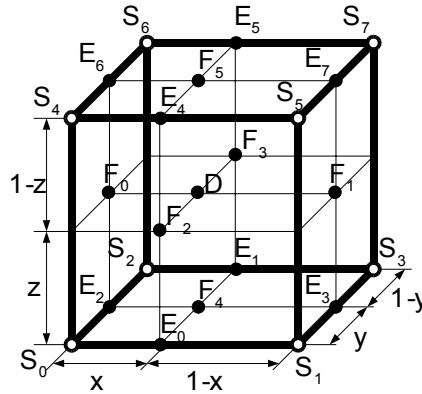


Abbildung 6.6: Interpolation und Gradientenschätzung nach Knittel

Eine der am einfachsten zu realisierenden Möglichkeiten der trilinearen Interpolation (Abbildung 6.6), bei der gleichzeitig die Gradienten abgeschätzt werden können, ist von Knittel in [7] veröffentlicht und wurde auch in dieser Arbeit verwendet.

Zuerst werden durch lineare Interpolation der acht Volumenelemente S_0 bis S_7 die Werte E_0 bis E_7 bestimmt. Aus E_0 bis E_7 werden wieder durch lineare Interpolation die Zwischenwerte F_0 bis F_5 berechnet:

$$\begin{aligned}
 E_0 &= S_1 x + S_0 (1-x) & F_0 &= E_6 z + E_2 (1-z) \\
 E_1 &= S_3 x + S_2 (1-x) & F_1 &= E_7 z + E_3 (1-z) \\
 E_2 &= S_2 y + S_0 (1-y) & F_2 &= E_4 z + E_0 (1-z) \\
 E_3 &= S_3 y + S_1 (1-y) & F_3 &= E_5 z + E_1 (1-z) \\
 E_4 &= S_5 x + S_4 (1-x) & F_4 &= E_3 x + E_2 (1-x) \\
 E_5 &= S_7 x + S_6 (1-y) & F_5 &= E_7 x + E_6 (1-x) \\
 E_6 &= S_6 y + S_4 (1-y) \\
 E_7 &= S_7 y + S_5 (1-y)
 \end{aligned}$$

Hierbei entsprechen x, y und z dem Abstand des Abtastpunktes zu den Gitterpunkten des Originalkoordinatensystems, welche der Position der Volumenelemente entsprechen. Die Abstände werden aus den Nachkommastellen der Abtastpunkt-Koordinaten gewonnen. Die gesuchten Werte für die Dichte D und die Komponenten des Gradientenvektors G_x, G_y, G_z werden wie folgt berechnet:

$$\begin{aligned}
 D &= F_{5z} + F_4 (1-z) \\
 G_x &= F_1 - F_0 \\
 G_y &= F_3 - F_2 \\
 G_z &= F_5 - F_4
 \end{aligned}$$

Diese Methode lässt sich sehr leicht realisieren und liefert mit acht parallelen Speichermodulen für jeden Zugriff einen interpolierten Wert und die zugehörigen Gradienten. Allerdings sind diese nicht sehr genau, da sie nur aus einer kleinen Umgebung gewonnen werden. In fast homogenen Bereichen haben sie deshalb eine nahezu willkürliche Orientierung [51]. Weiterhin entsprechen die Gradienten nicht genau der Position des Abtastpunktes, da sie nicht von allen Komponenten der Position des Abtastpunktes abhängig sind. So ist der Gradient in X-Richtung von der X-Komponente der Position des Abtastpunktes unabhängig und berücksichtigt mit der Interpolation der E- und F-Werte nur die Y- und Z-Position. Entsprechendes gilt auch für den Y- und Z-Gradienten. Nur mit größerem Hardware-Aufwand und vor allem dem Auslesen einer größeren Umgebung, lassen sich diese Nachteile, die zu einer schlechteren Bildqualität mit Artefakten führen, beseitigen.

Ebenfalls in [7] wird eine Berechnung der Gradienten über einen größeren Bereich vorgeschlagen. Hierbei entsprechen die Gradienten direkt dem Abtastpunkt und sind wesentlich genauer. Allerdings muss selbst mit einer Speicherarchitektur, die acht Volumenelemente parallel liefern kann, viermal pro Abtastpunkt auf den Speicher zugegriffen werden.

Bosma [52] zeigt eine Berechnung der Gradienten, die mit drei Zugriffen auskommt. Mit Hilfe von Look-Up-Tabellen werden vier Gradienten für eine Richtung über drei Stützstellen nichtlinear interpoliert. Diese vier werden wiederum linear auf den Abtastpunkt interpoliert. Durch die nichtlineare Interpolation werden harte Übergänge zwischen den Stützstellen geglättet und so Artefakte verringert.

Für ein System, das Bildraten für einen interaktive Umgang mit den Daten erzeugen soll, ist ein mehrmaliger Speicherzugriff pro Abtastpunkt nicht akzeptabel und höchstens als Option zu berücksichtigen.

Eine weitere Möglichkeit den Gradienten zu berechnen, ist nach der Interpolation und Post-Klassifikation. Dadurch wird automatisch, durch die Interpolation mit acht Elementen, ein größerer Volumenbereich für die Berechnung verwendet und die Rauschempfindlichkeit verringert. Allerdings sind bei der Realisierung zur Berechnung aller drei Komponenten des Gradienten große Zwischenspeicher notwendig, weshalb beim VIRIM-System (Kapitel 4.3) auf die Z-Komponente, die ein Zwischenspeichern einer ganzen Bildebene erfordert, verzichtet wird. Bei perspektivischer Projektion entstehen weitere Ungenauigkeiten, da kein orthogonales Koordinatensystem vorliegt [4]. Werden die Opazitäten durch die Post-Klassifikation nichtlinear verändert, können nur Gradienten, die nach der Post-Klassifikation berechnet wurden, korrekt sein (siehe Kapitel 5.2).

Wegen der großen Zwischenspeicher, die eine Realisierung in einer FPGA-basierten Hardware schwierig machen, wurde diese Variante, welche die qualitativ Beste ist, verworfen.

Unabhängig von der Art, wie die Gradienten berechnet werden, wird der Abtastwert nur aus den umgebenden acht Volumenelementen interpoliert.

Werden die acht Werte S_0 bis S_7 parallel von einer Speicherarchitektur (siehe Kapitel 7.5) geliefert, muss eine Zuordnung der acht Datenbusse D_0 bis D_7 zu den S-Werten erfolgen. Abbildung 6.7 zeigt einen Ausschnitt des Gitters mit den Volumenelementen und ihre Zuordnung zu den Datenbussen. Ein Volumenelement ist fest einem Datenbus zugeordnet und kann somit je nach Lage des Abtastpunktes jedem S-Wert entsprechen. Beispielsweise muss der Wert von S_0 , der der linken-unteren-vorderen Ecke entspricht, im Subwürfel 1 dem Datenbus D_0 entnommen werden. Wogegen S_1 , die rechte-untere-vordere Ecke, auf dem Datenbus D_1 liegt. Im Subwürfel 2 ist die Zuordnung aber gerade vertauscht, was erkennbar ist, wenn man sich den 2^3 -Subwürfel aus Abbildung 6.6 in Gedanken auf die Subwürfel in Abbildung 6.7 legt.

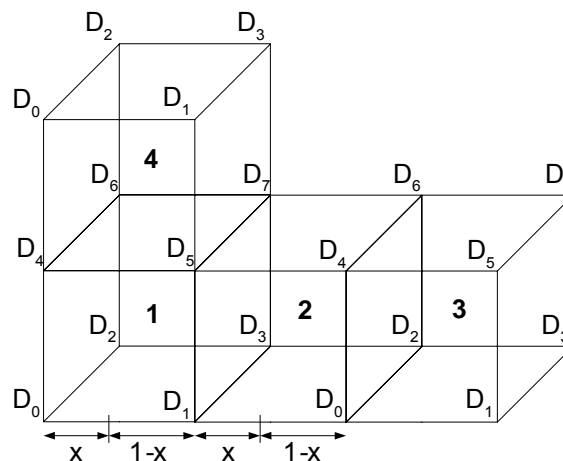


Abbildung 6.7: Zuordnung der Volumenelemente zu den Datenbussen

Eine mögliche Lösung zur Realisierung wären acht Multiplexer, welche die entsprechende Zuordnung der S-Werte zu den Datenbussen umsetzt. Weit weniger aufwändig sind allerdings drei Multiplexer für die Abstände x , y und z und ihren Differenzen $1-x$, $1-y$ und $1-z$. Das folgende Beispiel verdeutlicht, dass beide Vertauschungsmöglichkeiten zum gleichen Ergebnis führen:

Vertauschen der Datenbusse:

$$\text{in Subwürfel 1: } E_0 = S_{1x} + S_0(1-x) \rightarrow E_0 = D_{1x} + D_0(1-x)$$



$$\text{in Subwürfel 2: } E_0 = S_1 x + S_0(1-x) \rightarrow E_0 = D_0 x + D_1(1-x)$$

oder Vertauschen von x mit $1-x$:

$$\text{in Subwürfel 1: } E_0 = S_{1x} + S_0(1-x) \rightarrow E_0 = D_{1x} + D_0(1-x)$$



$$\text{in Subwürfel 2: } E_0 = S_1 x + S_0(1-x) \rightarrow E_0 = D_1(1-x) + D_0 x$$

Zudem werden bei den meisten Realisierungen 8 Bit für die x -, y - und z -Abstände ausreichen, während bei den S-Werten auch größere Wortbreiten realisiert werden. Wird die Zuordnung des Volumenelementes zum Datenbus über die drei niederwertigsten Bit der Rasterkoordinaten (x^0, y^0, z^0) kodiert^I, kann mit diesen auch das Vertauschen der x -, y - und z -Abstände gesteuert werden.

Die Umsetzung einer Interpolationsschaltung und Gradientenberechnung nach Knittel ist in Abbildung 6.8 gezeigt. Die Blöcke zur Berechnung der E- und F-Werte und der Dichte sind aus der Interpolationsschaltung aus Abbildung 6.3 aufgebaut; für die Gradientenberechnung sind noch drei zusätzliche Subtrahierer notwendig. Ein Block berechnet vorab die Abstände zu den Gitterpunkten $1-x$ und $1-y$ und gibt in Abhängigkeit von den niederwertigsten Bit der Gitterkoordinaten x^0 und y^0 , diese oder x und y an die Interpolation für die E-Werte weiter. Die Laufzeit der einzelnen Interpolationsstufen muss durch Verzögerungsglieder (τ) ausgeglichen werden, damit die Abstände der Gitterpunkte zu den richtigen Zwischenergebnissen passen.

^I vergleiche Kapitel 6.3 und Kapitel 7.5 Tabelle 8.2

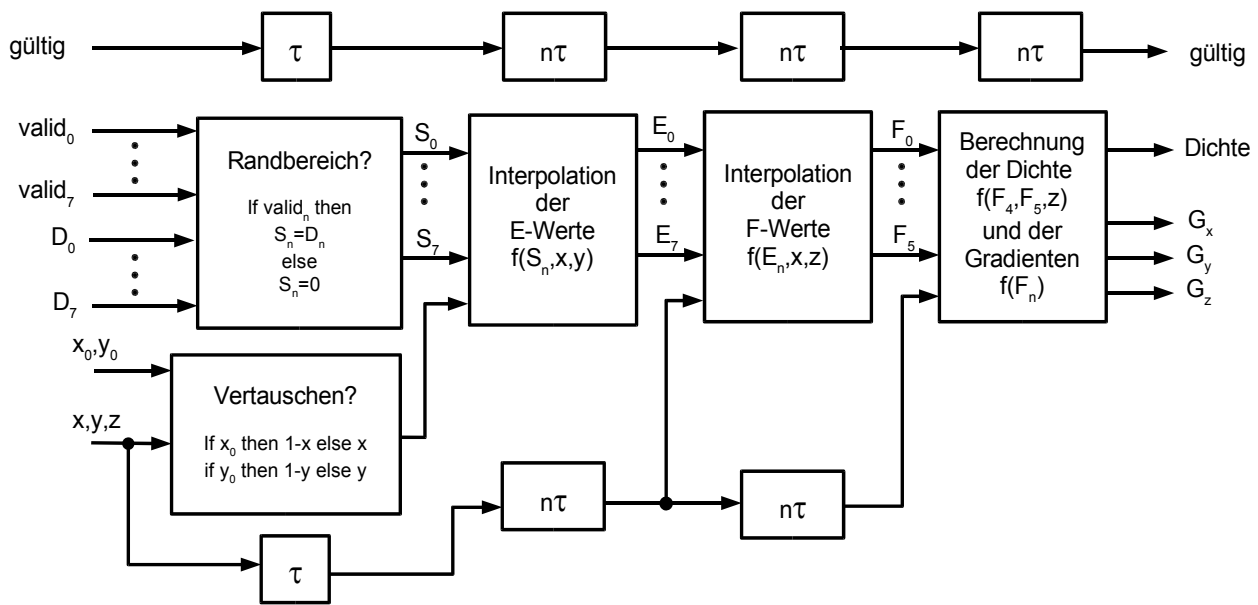


Abbildung 6.8: Gesamtsystem der Interpolation und Gradientenberechnung

Den Datenbussen D_0 bis D_7 ist ein Block zur Überprüfung nachgeschaltet, der S-Werte auf Null setzt, falls sie außerhalb des Speicherbereichs liegen. Dies wird an einem Überlauf bei der Adressberechnung festgestellt, wobei für jeden Datenbus ein *valid*-Signal generiert wird. Die Interpolation- und Gradienten-Pipeline arbeitet kontinuierlich, unabhängig, ob gültige Werte am Eingang anliegen oder nicht. Ein Flag (*gültig*) wird über alle Pipeline-Stufen verzögert, um am Ausgang anzuzeigen, wann die Dichte und Gradienten verwendet werden dürfen. Die Zahl n in Abbildung 6.8 entspricht der Anzahl der Verzögerungsglieder, die bei der Interpolationsschaltung verwendet wurden.

Die Schaltung ist fähig, in jedem Takt einen interpolierten Abtastwert mit zugehörigen Gradienten zu generieren.

6.6 Klassifizierung

Durch die Klassifizierung wird den Dichtewerten der Originaldaten ein Opazitätswert als Maß für die Undurchsichtigkeit zugeordnet. Da vom Benutzer jede beliebige Funktion eingestellt werden kann, bietet es sich an, dies über ein SRAM^I als Look-Up-Tabelle (LUT) zu realisieren.

Die Post-Klassifikation ist dabei etwas aufwändiger, da parallel alle 8 Datenbusse der

^I SRAM: **S**tatic **R**andom **A**ccess **M**emory, schneller flüchtiger Speicherbaustein. Statisch bezieht sich auf den Erhalt der Daten, die außer der Betriebsspannung keine weiteren Signale zur Auffrischung benötigen.

Speichermodule über SRAM-Bausteine geführt werden müssen. Die Datenbusse der SDRAM^I-Module werden dabei an die Adressleitungen der SRAM-Bausteine angeschlossen und deren Datenbusleitungen zur Interpolation und Gradientenberechnung weitergeführt. Vor der Visualisierung wird die Klassifizierung eingestellt, in dem an jeder möglichen SRAM-Adresse und somit jedem Datenwert der Originaldaten, ein beliebiger Wert eingespeichert wird. Über das SRAM kann hier auch eine Datenreduktion von z.B. 16 Bit der Originaldaten auf 8 Bit Opazitätswerte erfolgen. Dies reduziert die notwendige Eingangspinzahl und reduziert auch den Flächenaufwand für die Interpolationseinheit, vor allem bei einer FPGA-Zielhardware.

Zur Initialisierung der SRAM-Bausteine und des Volumenspeichers, muss deren Verbindung über Multiplexer trennbar und über einen Kontrollbuss ansprechbar sein. So müssen über die Datenleitungen der SDRAM-Bausteine die Volumendaten eingespeichert werden, und beim Einstellen der Klassifizierung die Adressleitungen der SRAM-Bausteine angesteuert werden. In Abbildung 6.9 ist eine Realisierungsmöglichkeit mit einem Transceiver-IC^{II} an einem einzelnen Speichermodule abgebildet. Das IC verbindet 2 Busse mit Treiberstufen, die wahlweise einzeln (*dir*-Eingang) oder beide (*sel*-Eingang) in den hochohmigen Zustand geschaltet werden können.

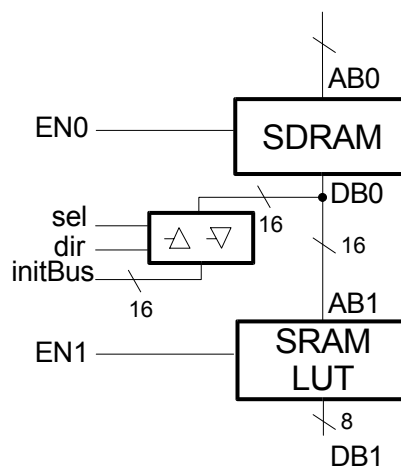


Abbildung 6.9: Multiplexer mit Transceiver-IC

Beim Einschreiben oder Auslesen der Volumendaten wird somit der *initBus* in der entsprechenden Richtung mit dem Datenbus (*DB0*) der SDRAMs verbunden und das SRAM-LUT über seinen Chip-Enable (*EN1*) abgeschaltet. Beim Konfigurieren der SRAM-LUT werden die SDRAMs mit *EN0* abgeschaltet und die Adressen auch über den *initBus* an den Adressbus (*AB1*) der SRAM-LUT

I SDRAM: Synchronous Dynamic Random Access Memory, Massenspeicher, der im Gegensatz zum SRAM eine regelmäßige Auffrischung der Daten benötigt, genauere Erläuterung in Kapitel 7.3.2

II z.B. SN74ALB16245 von Texas Instruments

angelegt. Während der Visualisierung, bleibt das Transceiver-IC vollständig passiv (hochohmig).

Wie Hardware-Testaufbauten zeigten, sind asynchrone SRAMs nur mit aufwändigen Timing-Überprüfungen eines FPGA-Designs verwendbar. Bei stark ausgenutzten FPGAs war das notwendige Zeitverhalten kaum einzuhalten. Es kommt dann zu sporadischen Schreib- und Lesefehlern. Deshalb sind synchrone SRAMs, die mit einer taktsynchronen Schnittstelle erweitert wurden, zu bevorzugen.

Sehr einfach ist die Post-Klassifikation als SRAM-LUT im FPGA zu realisieren, da nach der Interpolation nur noch ein einzelner Datenbus vorliegt. Dadurch reicht bei einer 8-Bit Auflösung ein 256 Byte internes, synchrones SRAM aus.

6.7 Shading

In Kapitel 5.3 wurden die Beleuchtungsmodelle und deren Berechnung erläutert. Während der Visualisierung muss die Reflexion zum Beobachter für jeden Abtastpunkt neu berechnet werden, wobei vor allem der spiegelnde Anteil einigen Rechenaufwand erfordert. Anstatt ein entsprechend großes Rechenwerk aufzubauen, gibt es auch die Möglichkeit, die Berechnungen vorher durchzuführen und die Ergebnisse in einer Look-Up-Tabelle abzulegen [15].

Während eines Bildes bleibt die Beobachter- und Lichtrichtung konstant. Dadurch ergibt sich bei der Reflexionsberechnung nur noch eine Abhängigkeit von der Ausrichtung der Oberfläche zum Beobachter. Man kann dadurch für gleichmäßig verteilte Punkte auf der Oberfläche einer Kugel die Reflexionen vorberechnen und somit alle möglichen Ausrichtungen einer Oberfläche zum Beobachter abdecken. Um die Zahl der Punkte klein zu halten, kann man sie als Stützstellen für eine lineare Interpolation verwenden und so auch ohne große Fehler Zwischenwerte erhalten.

Die Stützstellen werden in einer Tabelle abgelegt. Die Tabelle wird für einen Beobachtungsvektor und beliebig viele Beleuchtungsquellen vorberechnet und muss bei Änderung einer dieser Vektoren neu berechnet werden. Da man mit einer kleinen 1.5 KByte großen Tabelle auskommt, kann sie bei Rotation des Objektes für jedes Bild neu berechnet werden. Innerhalb eines Bildes gilt die Tabelle exakt allerdings nur für parallele Beobachter- und Beleuchtungsvektoren, da sich bei Perspektive der Beobachterwinkel über das Bild ändert.

Es muss noch geklärt werden, wie eine Kugeloberfläche am besten auf eine Tabelle abgebildet werden kann und wie diese adressiert wird.

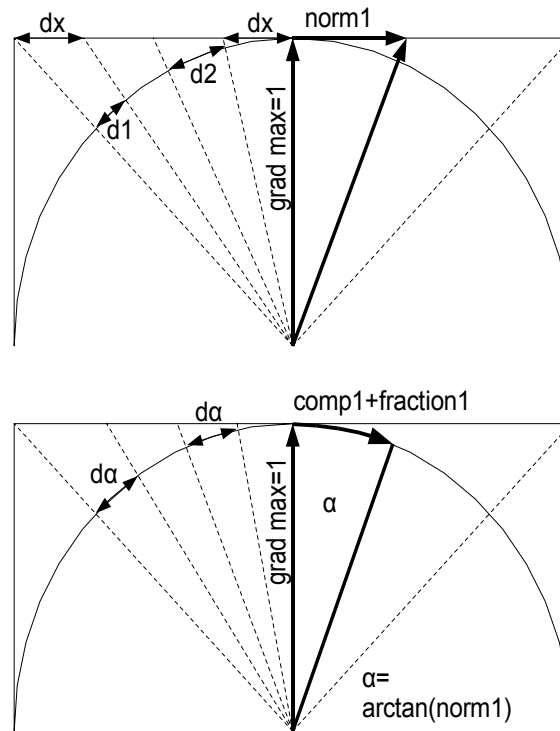


Abbildung 6.10: Adressierung der Reflexionstabelle mit und ohne Distorsions-Kompensation.

Die Ausrichtung und das Vorhandensein einer Oberfläche wird durch den Gradienten am Abtastpunkt bestimmt. Wird nur die Richtung des Gradienten berücksichtigt und dieser auf 1 normiert, entspricht dies der Oberflächennormalen, die auf einen beliebigen Punkt auf die Oberfläche einer Einheitskugel zeigt. Der Abtastpunkt liegt in deren Zentrum. Da eine Kugel schwer auf eine Tabelle abzubilden ist, wird um die Kugel ein Würfel gelegt. Um nun mit dem Gradienten einen Punkt auf der Würfeloberfläche auszuwählen, werden die drei Komponenten des Gradienten durch die größte Komponente dividiert. Dieser Vorgang wird als Normalisierung bezeichnet und hat als Ergebnis die beiden kleineren normierten Komponenten $norm_1$ und $norm_2$. Die größte Komponente des Gradienten ($gradmax$) hat danach die Länge 1 und gibt mit dem Vorzeichen die Würfelseite (+/-x, +/-y oder +/-z) an, auf die der Gradient zeigt. Die beiden kleineren Komponenten geben dann einen Punkt auf der ausgewählten Würfelfläche an.

Die Adressierung dieser Tabelle ist sehr einfach, hat aber den Nachteil, dass die Winkelauflösung der Gradienten bei 0° kleiner ist als bei 45° . Dieser Effekt, den man *Distorsion* nennt [15], ist in Abbildung 6.10 oben an einer unterschiedlichen Länge von $d1$ und $d2$ zu erkennen, während dx auf der Würfeloberfläche gleich bleibt. Auf die Bildqualität wirkt sich dies nur leicht aus, kann aber

einfach durch zwei weitere Look-Up-Tabellen vermieden werden, indem man die beiden kleineren, normierten Gradientenkomponenten ($norm_{1,2}$) über eine \arctan -Funktion in einen Winkel ($\alpha_{1,2}$) umwandelt. Die Tabelle wird dann über die Winkelgröße kodiert. Die Kodierung von 0° - 45° und ein Vorzeichen ist hierbei ausreichend. In Abbildung 6.10 unten wird dies durch die gleichen Längen von $d\alpha$ veranschaulicht.

Wobei $\tan(\alpha_{1,2}) = \frac{norm_{1,2}}{grad_{max}} = \frac{norm_{1,2}}{1} = norm_{1,2}$ und somit $\alpha_{1,2} = \arctan(norm_{1,2})$ ist.

Abbildung 6.11 zeigt das Blockdiagramm der kompletten Schaltung aus [15].

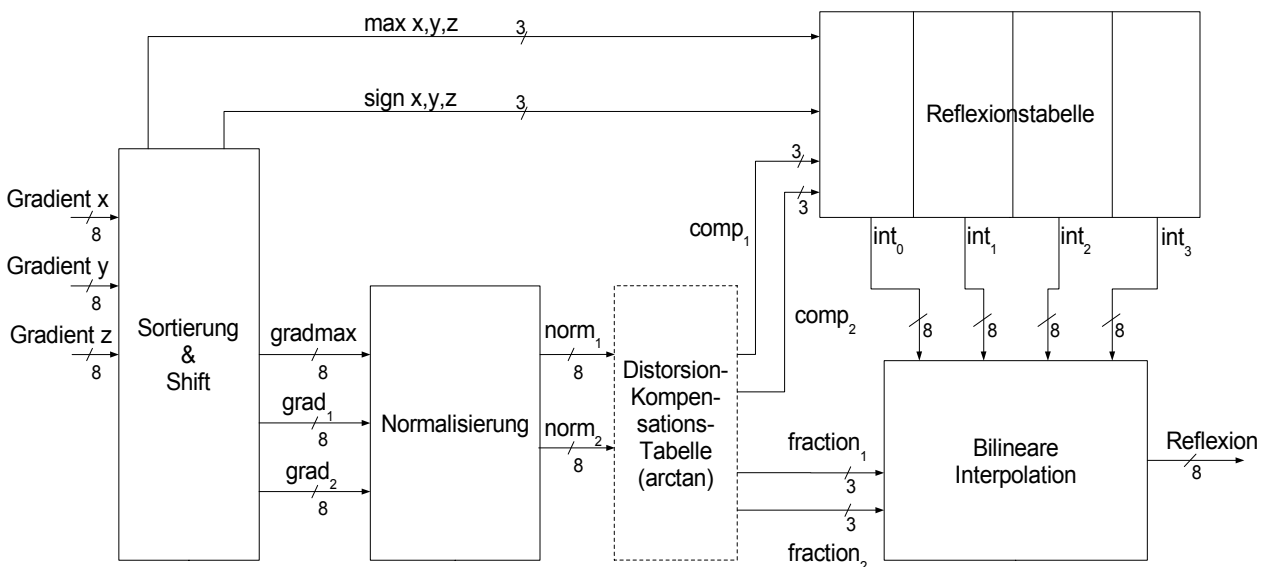


Abbildung 6.11: Blockschaltbild der Reflexionsberechnung

Zuerst wird im Block *Sortierung & Shift* aus den drei Komponenten des Gradienten die größte herausgesucht und eventuell führende Nullen gleichmäßig über alle Komponenten herausgeschoben, um die größte Genauigkeit zu erreichen. Die Vorzeichenbit ($sign\ x,y,z$) und je eine Leitung zur Signalisierung der größten Komponente ($max\ x,y,z$) werden direkt zur Adressierung der Würfelseite aus der Tabelle verwendet^I. Zur weiteren Berechnung werden nur die Beträge von $grad_1$, $grad_2$ und $grad_{max}$ weitergeben.

Bei der Normalisierung werden die kleineren Komponentenanteile durch den größten dividiert. Mit den höherwertigen Bit der Ergebnisse ($norm_1$, $norm_2$) könnte ohne Distorsion-Kompensationstabelle bereits die Reflexionstabelle adressiert werden. Die niederwertigen Bit würden dann zur

^I So wird es in [15] beschrieben, da allerdings nicht alle Kombinationen von $max\ x,y,z$ und den drei Vorzeichenbit vorkommen, können die sechs Bit auf drei Bit zur binären Kodierung der sechs Würfelseiten reduziert werden.

Zwischenwertberechnung bei der Interpolation verwendet werden (siehe auch Abbildung 6.10 oben).

Mit Distorsion-Kompensation werden $norm_1$ und $norm_2$ getrennt über eine Look-Up-Tabelle mit der arctan-Funktion in einen Winkel ($comp_{1,2} + fraction_{1,2}$) umgewandelt.

Aus der Reflexionstabelle werden parallel mit $comp_1$ und $comp_2$ die vier umliegenden Stützwerte ($int_{0..3}$) der Würfel- bzw. Kugeloberfläche an die *bilineare Interpolation* weitergegeben. Der Abstand zu den Stützstellen in zwei Richtungen wird durch $fraction_1$ und $fraction_2$ angegeben. Es wird wieder die gleiche Interpolationsschaltung, wie in Abbildung 6.3 gezeigt, verwendet.

Die Größe der Reflexionstabelle berechnet sich aus 6 Würfel- bzw. Kugelabschnitten, ausgewählt durch $max\ x,y,z$ und $sign\ x,y,z$, multipliziert mit 2^6 Adressen durch $comp_{1,2}$ mit je 4 Stützstellen zu 8 Bit und ergibt 1,5kByte.

Die Tabelle wird vorberechnet, indem für jede Adresse die zugehörige Oberflächennormale bestimmt wird. Mit vorgegebenem Beobachter- und Beleuchtungsvektor kann nach Phong die Reflexion bestimmt werden. Bei der Adressierung der Tabelle werden parallel 4 Reflexionen auf der Kugeloberfläche ausgegeben. Die umliegenden Stützstellen ergeben sich durch Erhöhen der Adressen $comp_1$ und $comp_2$ (siehe Tabelle 6.1) und können zusammen als 32-Bit Wort in ein entsprechendes RAM geschrieben werden.

int_0	$comp_1$	$comp_2$
int_1	$comp_1+1$	$comp_2$
int_2	$comp_1$	$comp_2+1$
int_3	$comp_1+1$	$comp_2+1$

Tabelle 6.1 Zuordnung der Stützstellen zur Adressierung

Dadurch ist jede Stützstelle zwar viermal in der Tabelle enthalten, aber aufgrund der geringen Ausmaße der Tabelle ist das ohne Probleme möglich.

7 Grundprinzip der Architektur für Optimierungen in Hardware

7.1 Pipelining – Hazard

Abbildung 7.1 zeigt den schematischen Ablauf einer Standard Ray-Casting-Pipeline, die einen Strahl durch das Volumenobjekt sendet und das zugehörige Pixel im Ergebnisbild berechnet.

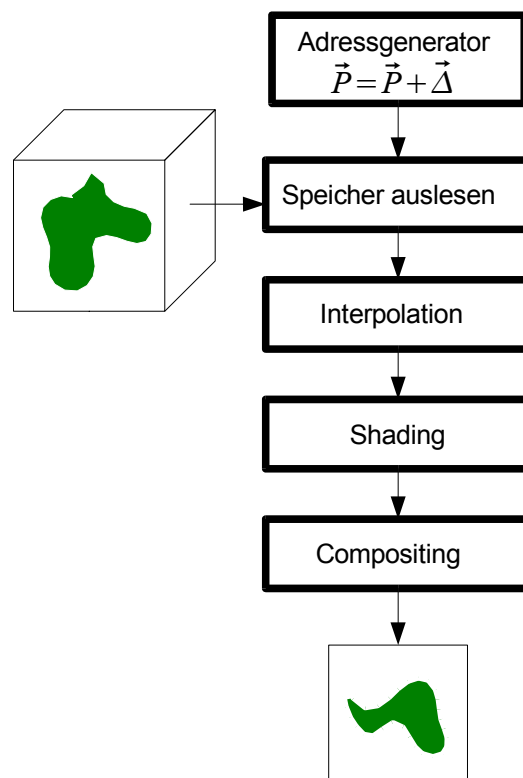


Abbildung 7.1: Standard Ray-Casting-Pipeline ohne Optimierungen

Ausgangspunkt ist die Adressgenerierung. Sie startet an einem Punkt \vec{P} , der in der Projektionsebene auf einem Bildpunkt liegt. In der nächsten Stufe werden die dem Punkt \vec{P} entsprechenden Daten aus dem Speicher gelesen und schrittweise an die darauf folgenden Stufen Interpolation, Shading und Compositing weitergeleitet. Während die Daten aus dem Speicher gelesen werden, berechnet die Adressgenerierung den nächsten Punkt durch Addition eines Differenzvektors $\vec{\Delta}$, der die Beobachtungsrichtung und dessen Betrag $|\vec{\Delta}|$ den Abtastabstand angibt. In jedem Arbeitstakt der Pipeline werden neue Daten angefordert und durch die Pipeline geschoben. Nachdem die Pipeline komplett gefüllt ist, kann in jedem Takt ein neuer Wert im Compositing berechnet werden. Ist der Strahl komplett abgearbeitet, wird das Ergebnis-Pixel

gespeichert und die Adressgenerierung für den nächsten Strahl initialisiert. Bei diesem sequentiellen Aufbau ohne Rückkopplung im Datenfluss könnte die Initialisierung der Adressgenerierung erfolgen, während noch Daten des vorherigen Strahls bearbeitet werden. Die Pipeline kann dadurch ununterbrochen ohne Leerlauf arbeiten.

Wie sieht es aber aus, wenn der Strahl nicht völlig abgearbeitet, sondern mit Hilfe algorithmischer Optimierungen, nur die notwendigen Teile des Strahls berechnet werden sollen?

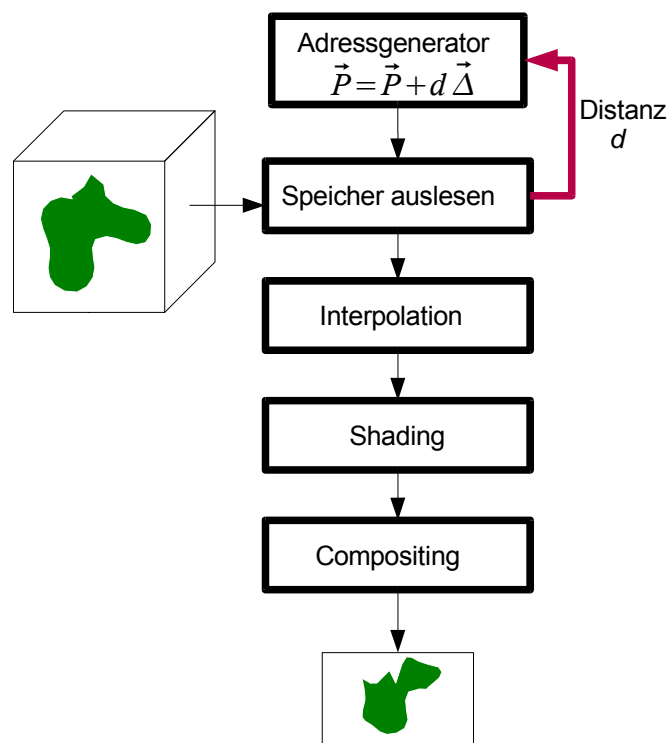


Abbildung 7.2: Ray-Casting Pipeline mit Auslesen der Sprungdistanz für Space-Leaping

Abbildung 7.2 soll das Problem am Beispiel von Space-Leaping als Optimierungstechnik verdeutlichen.

Bei Space-Leaping werden für leere Volumenbereiche die minimalen Sprungdistanzen zu den nächsten gefüllten Bereichen gespeichert (siehe Kapitel 3.2). Beim Visualisieren wird mit jedem Abtastpunkt zusätzlich ein Distanzwert d ausgelesen, der in der Adressgenerierung, um den nächsten Abtastpunkt zu berechnen, mit dem Differenzvektor $\vec{\Delta}$ multipliziert wird. Bei diesem Aufbau befindet sich im Unterschied zur Standard Ray-Casting-Pipeline eine Rückkopplung im Datenfluss. Ohne zusätzliche Änderung würden Wartezeiten in der Pipeline entstehen, da die Adressgenerierung nicht mehr parallel zum Datenauslesen den nächsten Abtastpunkt berechnen kann, sondern warten muss bis der nächste Distanzwert zur Verfügung steht.

Aktuelle Speicherbausteine besitzen ebenfalls eine Pipelinearchitektur. Dadurch werden nach dem Anlegen eines Lese-Befehls erst drei Takte später die Daten ausgegeben. Zusätzlich wird die Berechnung des neuen Abtastpunktes ebenfalls zwei Takte benötigen, für Multiplikation des Sprungwertes mit dem Distanzvektor $d\vec{\Delta}$ und für die Addition zum vorherigen Abtastpunkt $\vec{P} + d\vec{\Delta}$. Dadurch könnte erst nach den fünf Takten wieder ein neuer Lesebefehl gegeben werden.

Kommt es zu einer Unterbrechung einer kontinuierlichen Berechnung spricht man von einem *Hazard* (engl. Unglücksfall). Im Fall von Space-Leaping muss die Pipeline Wartezyklen einschieben, da auf fehlende Daten gewartet werden muss. Es handelt sich hierbei um einen *Daten-Hazard*.

Bei Early-Ray-Termination wird im Verlauf der Verarbeitung festgestellt, dass die weitere Berechnung abgebrochen werden kann. Und zwar, wenn durch die Absorption, die Lichtintensität unter einen Schwellwert gefallen ist, so dass kein nennenswerter Beitrag zum Bild hinzukommt. Es werden alle, momentan in der Pipeline befindlichen Daten ungültig und die Verarbeitung muss zu einem neuen Strahl "verzweigen", weshalb man von *Branch-Hazard* spricht.

Beide Fälle führen dazu, dass die Pipeline nicht kontinuierlich arbeiten kann und so die Einsparungen durch die algorithmischen Optimierungen, wieder hinfällig werden. Es gibt aber Lösungsmöglichkeiten für dieses Problem.

7.2 Multithreading zur Steigerung der Auslastung

Abbildung 7.3 oben zeigt in vereinfachter Darstellung die Auslastung der Pipeline über der Zeit. Die weißen Kästchen zeigen wann die drei Arbeitsschritte Speicherauslesen, Multiplikation und Addition inaktiv sind. In diesem Beispiel wurde für das Auslesen des Speichers nur ein Takt veranschlagt. Man kann gut erkennen, dass die Pipeline in diesem Beispiel nur zu einem Drittel ausgelastet ist.

In der Mikroprozessortechnik gibt es ein ähnliches Problem der Auslastung, wenn Programme beispielsweise auf Ein- oder Ausgabe warten, muss der Prozessor Wartezyklen einschieben. Lösungen hierfür ergaben sich durch Multiprozess-Betriebssysteme und Prozessoren, die einen schnellen Prozesswechsel ermöglichen, falls der aktive Prozess in einen Wartezustand geht. Gleichzeitig laufende Prozesse, die mit den gleichen Daten arbeiten, werden Thread (engl. Faden) genannt. Multithreading nennt sich das Prinzip des schnellen Wechsels zwischen mehreren, parallel laufenden Threads.

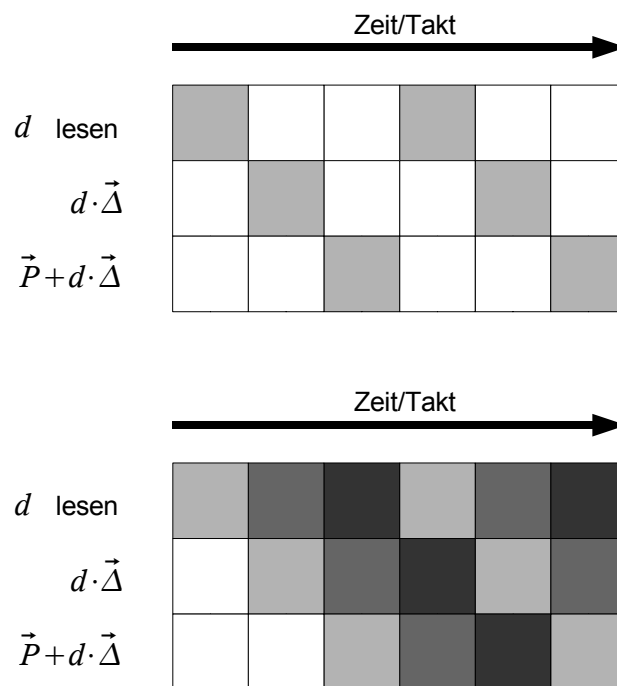


Abbildung 7.3: Auslastung der Pipeline

Die Idee war nun, dieses Prinzip auf eine Ray-Casting-Pipeline anzuwenden und dadurch die volle Auslastung zu erreichen, auch wenn durch algorithmische Optimierungen kein vorhersehbarer Ablauf einzelner Sehstrahlen möglich ist.

Abbildung 7.3 unten zeigt das Prinzip. Gleiche Graustufen gehören zum selben Sehstrahl. In den Leerlaufphasen der Pipeline werden weitere Sehstrahlen eingeschoben und so die Wartezeiten überbrückt und wieder eine volle Auslastung erreicht.

Es stellt sich die Frage, wie der Strahlwechsel am besten vorzunehmen ist.

7.3 Ausnutzung der Datenkohärenz

Erste Simulationen unter der Berücksichtigung des Zeitverhalten realer Speicher zeigten, dass ein einfacher Wechsel der Sehstrahlen reihum nicht den gewünschten Erfolg bringt.

Die Speicher sind in Blöcke oder Zeilen aufgeteilt und haben die Eigenschaft, Daten innerhalb einer Zeile sehr schnell zu liefern. Bei Zugriffen auf eine andere Zeile allerdings entstehen erhebliche Wartezeiten. Um nicht alle Vorteile der Optimierungstechniken an diese Wartezeiten zu verlieren, ist eine sehr genaue Abstimmung der Strahlwechsel mit der Speicherarchitektur notwendig.

Man muss dafür sorgen, dass die vom Algorithmus nacheinander benötigten Daten möglichst zusammenhängend im Speicher liegen und, falls möglich, der Algorithmus benachbarte Daten

bevorzugt bearbeitet. Beides zusammen wird hier als Ausnutzung der Datenkohärenz verstanden.

7.3.1 Vergleich der Speicherbausteine

Für das Volume-Rendering-Verfahren müssen große Datenmengen gespeichert und mit hohem Durchsatz ausgelesen werden. Ein Datensatz mit 1024^3 Volumenelementen bei 16 Bit Auflösung benötigt 2 Gigabyte Speicherplatz. Statische Speicher sind für diese Datenmengen nicht brauchbar und können höchstens als Zwischenspeicher verwendet werden. Daher kommen nur kostengünstige dynamische Speicherbausteine in Frage.

Zur Zeit auf dem Markt erhältlich sind im wesentlichen SDR^I- und DDR^{II}-SDRAMs^{III}. Beide sind in der Ansteuerung im wesentlichen gleich. Der einzige Unterschied besteht im Block-Lesemodus, wobei DDR-SDRAM mit jeder Taktflanke – steigende und fallende – neue Daten aus geben, anstatt nur mit der steigenden Flanke.

Werden die Speicherbausteine ohne den Block-Lesemodus verwendet, also wahlfrei auf einzelne Adressen zugegriffen, bieten die DDR-SDRAMs keinen Geschwindigkeitsvorteil gegenüber den SDR-SDRAMs.

Bei DRDRAM^{IV} handelt es sich um einen weiteren Massenspeicher, der von der Firma Rambus lizenziert wird und einen sehr hohen Datendurchsatz ermöglicht. Die Bausteine benötigen einen speziellen Rambus Speicher Controller^V (RAC), der nur als Block für integrierte Schaltungen angeboten wird und dadurch nur für Produkte mit hohen Stückzahlen verwendet werden kann.

Wie bei den DDR-SDRAMs wird der volle Datendurchsatz der DRDRAMs nur beim Auslesen großer Blöcke erreicht.

Will man auf eine zusätzliche Cache-Architektur als Zwischenspeicher verzichten und die Daten wahlfrei aus dem Speicher lesen, sind SDR-SDRAMs am besten geeignet.

Die im Folgenden beschriebene Speicherarchitektur nutzt die interne Struktur der SDRAM-Bausteine als Zwischenspeicher, weshalb zum Verständnis eine Erläuterung ihrer Funktionsweise notwendig ist.

7.3.2 Single Data Rate SDRAM

Die Technologie von SDR SDRAM-Bausteinen basiert auf herkömmlichen dynamischen

I SDR: Single Data Rate; 125MB/s Datendurchsatz

II DDR: Double Data Rate; 200MB/s Datendurchsatz

III SDRAM: Synchronous Dynamic Random Access Memory

IV DRDRAM: Direct Rambus Dynamic Random Access Memory; pro Kanal 800MB/s Datendurchsatz, 4 Kanäle möglich

V RAC: Rambus ASIC Cell

Speicherbausteinen (DRAM), die durch eine Steuerlogik und Ein- und Ausgangspuffer um ein taktsynchrones Interface erweitert wurden. Sie werden in verschiedenen Konfigurationen angeboten, die sich in der Gesamtkapazität und der Datenbusbreite unterscheiden. Üblich sind momentan 4,8 und 16 Bit breite Datenbusse und Gesamtkapazitäten von 128, 256 und 512 Megabit. Alle Bauformen besitzen vier interne, unabhängige DRAM-Speicherbänke.

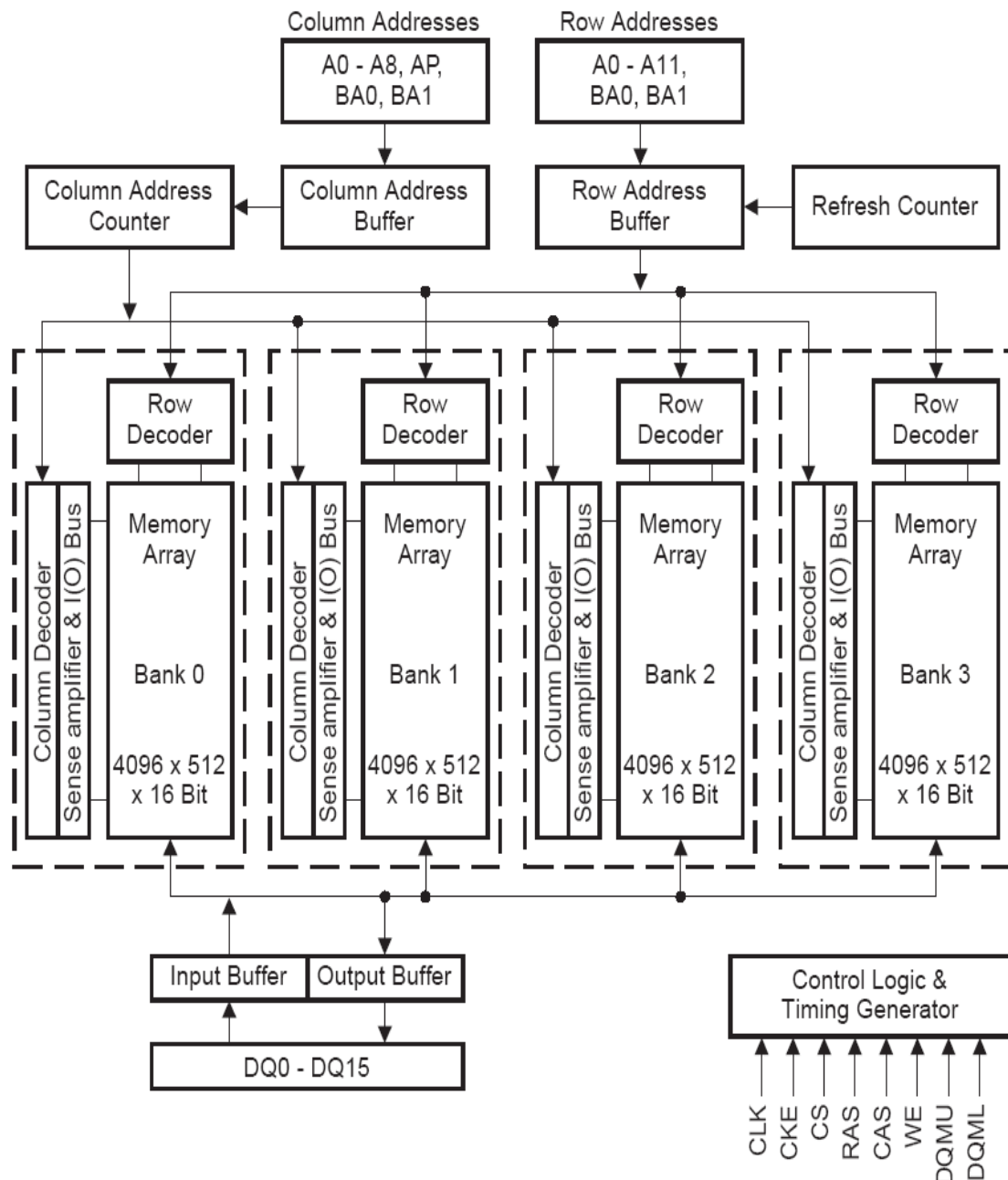


Abbildung 7.4: Blockdiagramm SDRAM Infineon HB39S128CT 4 Bänke, 4096 Zeilen, 512 Spalten mit 16 Bit Datenworten (128 MBit)

Abbildung 7.4 zeigt das Blockdiagramm eines SDRAM-Bausteins organisiert in vier Bänken mit je 4096 Zeilen (Row), 512 Spalten (Column) mit 16 Bit Datenbusbreite, was einer Gesamtkapazität

von 128 Megabit entspricht.

Die DRAM-Bänke (Memory Array) bestehen aus Bit-Zellen, die in Spalten und Zeilen organisiert sind. Jede Bit-Zelle wird durch einen Kondensator realisiert, der durch einen Transistor angesteuert werden kann. Der Zeilendekoder wählt anhand der zuvor in den Baustein eingelesenen Zeilenadresse die entsprechende Zeile aus, indem er die dazu notwendigen Transistoren aktiviert.

Die Ladung der Kondensatoren ist zu gering, als dass man sie direkt auslesen könnte; deshalb sind die DRAM-Bänke wiederum in zwei symmetrische Hälften (hier 2 x 2048 Zeilen) unterteilt, wobei immer nur in einer Hälfte eine Zeile ausgewählt wird. Komparatoren (Sense Amplifier) detektieren dann, für jede ausgewählte Bit-Zelle getrennt, den geringen Ladungsunterschied zwischen der ausgewählten und der nicht ausgewählten Bit-Zelle in den beiden Hälften. Zuvor müssen allerdings die Verbindungsleitungen zu den Transistoren auf die halbe Versorgungsspannung aufgeladen werden (Precharge) [36]. Fällt beim Detektieren (Activate), also beim Durchschalten eines Transistors, die Spannung etwas gegenüber der nicht angesteuerten Leitung ab, wird vom Komparator eine Null erkannt, steigt sie etwas an, wird eine Eins ausgegeben.

Die Komparatoren detektieren die Ladungen für eine komplette Spalte; hier zum Beispiel 512 Spalten mal 16 Bit entspricht 8192 Komparatoren^I.

Mit der Spaltenadresse werden über den Spaltendekoder 16 Ausgänge der Komparatoren ausgewählt und an die Ausgangspuffer übertragen. Eine Änderung der Spaltenadresse bei darauf folgenden Lesezugriffen stellt nur den Spaltendekoder neu ein, benötigt aber keine neue Messung durch die Komparatoren.

Beim Schreiben spielen die Komparatoren keine Rolle, hier werden die Kondensatoren direkt über die aktivierten Transistoren an den Eingangspuffern ge- oder entladen.

Durch die Selbstentladung der Kondensatoren muss alle 64 Millisekunden ein Refresh-Zyklus eingelegt werden, der beim Durchschalten aller Zeilenadressen die Kondensatoren wieder mit Ladung versorgt oder entlädt.

Mit diesem kurzen Einblick in die Arbeitsweise von DRAM lässt sich die Befehlsabfolge zur Ansteuerung von SDRAM-Bausteinen besser verstehen.

Alle Befehle an den SDRAM-Baustein werden zur steigenden Taktflanke, über die in Abbildung 7.4 an die *Control Logic* angeschlossenen Signale, eingelesen. Die wichtigsten davon sind CAS (Column Access Strobe), RAS (Row Access Strobe) und WE (Write Enable). Über die Signale BA0

^I Durch ein Flip-Flop im Komparator wird die Verbindungsleitung komplett auf die Versorgungsspannung oder Masse geschaltet und somit der Kondensator aufgefrischt. Dieser Vorgang wird auch bei Refresh ausgenutzt, wobei alle Zeilen einmal angesprochen werden und dabei die Kondensatoren der kompletten Spalte jeweils ge- oder entladen werden.

und BA1 (Bank Address) wird der Befehl einer der vier DRAM-Bänke zugeordnet.

Um ein Datenwort aus einem SDRAM zu lesen, müssen als erstes, für die entsprechende Bank in der sich das Datenwort befindet, die Verbindungsleitungen zu den Transistoren mit dem PRECHARGE-Kommando vorgeladen werden. Danach folgt das Einspeichern der Zeilenadresse und die Detektion durch die Komparatoren mit dem ACTIVE-Kommando. Dann kann in jedem weiteren Takt ein Datenwort aus dieser Bank mit der gleichen Zeilenadresse gelesen werden (READ-Kommando), ohne dass das PRECHARGE- oder ACTIVE-Kommando neu ausgeführt werden muss^I. Durch die Pipelinestruktur der SDRAM-Steuerlogik liegen die Datenwörter um drei Takte verzögert am Datenbus an.

Sind alle vier Bänke aktiviert, kann somit in jedem Takt wahlfrei von jeder Bank ein Datenwort gelesen werden.

Bei mit 100MHz getakteten SDRAM-Bausteinen benötigt das PRECHARGE-Kommando je nach Baustein 20-30ns, oder 3 Takte. Das ACTIVE-Kommando benötigt 50-60ns, also 6 Takte, und READ 1 Takt. Bei Änderung der Zeilenadresse müssen die Befehle PRECHARGE, ACTIVE und READ der Reihe nach angelegt werden, wodurch insgesamt 10 Takte benötigt werden, statt einem Takt, wenn man auf der gleichen Zeilenadresse arbeitet^{II}.

Man kann die aktivierten Bänke der SDRAM-Speicher auch als schnellen Zwischenspeicher (Cache) betrachten. Der in Abbildung 7.4 gezeigte Baustein besitzt beispielsweise 4 KByte (4x512x16Bit) Zwischenspeicher, der wahlfrei und ohne Verzögerung adressierbar ist. Gerade für ein FPGA-System, das über beschränkte SRAM-Ressourcen verfügt, ist es sinnvoll, die Zwischenspeicher auszunutzen.

In gleicher Weise können auch DDR-SDRAM-Bausteine verwendet werden, nur kann durch den wahlfreien Zugriff deren Geschwindigkeitsvorteil nicht ausgenutzt werden.

7.3.3 Ausnutzung aktiver SDRAM-Bänke als Zwischenspeicher

Die ersten Simulationen des Multithreading-Ansatzes (siehe Kapitel 7.2) zeigten, dass der erwartete Geschwindigkeitsgewinn durch ständiges Wechseln der Zeilenadresse im SDRAM-Speicher verloren geht.

Würde der Volumendatensatz sequentiell im Speicher abgelegt werden, also mit aufsteigender Speicheradresse erst die X, dann Y und dann die Z-Komponenten des Datenwürfels, gäbe es bei

I Es gibt eine Obergrenze von etwa 100us, dann muss PRECHARGE und ACTIVE wieder neu ausgeführt werden, da wegen Leckströmen die Kondensatoren neu geladen werden müssen.

II Ohne die 3 Takte Verzögerung zum Datenbus

jedem Abtastpunkt mehrere Zeilenwechsel, da die Volumendaten einer Ebene zeilenweise in den SDRAM-Zeilen liegen. Ein Abtastpunkt mit seinem 2^3 -Subwürfel benötigt aber Daten aus zwei Zeilen und zwei Ebenen (siehe Simulationsergebnis in Kapitel 9.1.6).

Um Richtungsunabhängigkeit und bestmögliche Datenkohärenz zu erreichen, müssen die Daten deshalb würfelförmig in einer SDRAM-Zeile abgelegt werden; bei 512 Volumenelementen pro Zeile entspricht dies einem 8^3 Würfel.

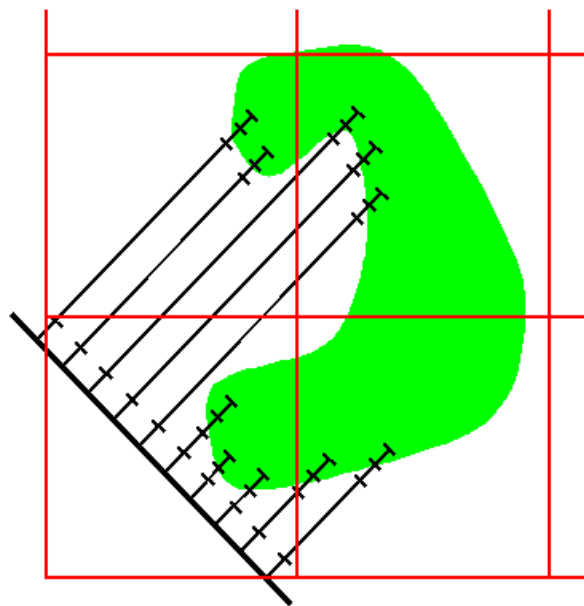


Abbildung 7.5: Sprung von Strahlen in unterschiedliche SDRAM-Zeilen bei Space-Leaping und würfelförmiger Einspeicherung

Weiterhin kann man beim Space-Leaping-Verfahren beobachten, dass häufig Strahlen mit ihrem nächsten Abtastpunkt in neue SDRAM-Zeilen springen, während andere Strahlen noch in der alten Zeile Daten benötigen. Abbildung 7.5 soll dies verdeutlichen. Die Quadrate stellen zusammengehörige Daten einer SDRAM-Zeile dar. Würde man die Strahlen reihum bearbeiten, müssten alle vier Bereiche regelmäßig neu aufgerufen werden.

Wesentlich besser wäre eine Strategie, welche die Sprungweite berücksichtigt, und nahe gelegene Abtastpunkte zuerst bearbeitet. Dies könnte zum Beispiel mit einer Sortierung nach dem Abstand zur Projektionsebene erreicht werden (siehe hierzu Simulationsergebnis in Kapitel 9.1.3).

Die Strahlen müssen nicht nur in der Länge zusammengehalten werden, sondern auch in Höhe und Breite, da es sich um ein räumliches Problem handelt. Deswegen sollten die Strahlen als Bündel die Bildpunkte eines kleineren, quadratischen Teilbildes bestimmen und beispielsweise nicht in einer Zeile liegen (Simulationsergebnis Kapitel 9.1.2).

Ein SDRAM-Baustein kann vier Zeilen parallel aktiv halten. Dies kann man ausnutzen, indem benachbarte Volumenbereiche in unterschiedlichen SDRAM-Bänken untergebracht werden. Dadurch kommt es selbst, wenn ein Strahl über eine Nachbargrenze kommt, nicht sofort zu einer Verdrängung der Zeile.

Mit dieser Kombination aus Sortierung der Strahlen und würfelförmigem Aufteilen des Volumens auf SDRAM-Zeilen und Bänke ist es möglich, die SDRAM-Zeilen als Zwischenspeicher zu verwenden. Dadurch bleiben die Zeilen solange aktiv, bis alle Strahlen den Datenbereich verlassen haben und im optimalen Fall nur einmal pro Strahlbündel geladen werden.

Im Folgenden wird zuerst auf die Realisierung der Sortierfunktion eingegangen, bevor anschließend eine ideale Speicher-Architektur vorgestellt wird.

7.4 Implementierung von Multithreading und Sortierfunktion

7.4.1 Erweiterung des Adressgenerators

Änderungen an der in Abbildung 7.1 gezeigten Ray-Casting-Pipeline, für Multithreading und Sortieren der Sehstrahlen, sind hauptsächlich am Adressgenerator notwendig. Statt eines Strahls, müssen die Parameter mehrerer Strahlen verwaltet und nach Abstand zur Projektionsebene sortiert werden. Abbildung 7.6 zeigt das Prinzip. Die Start- und Distanzvektoren aller Strahlen werden bei der Initialisierung in einem Parameterspeicher abgelegt.

Das Herz des erweiterten Adressgenerators besteht aus der Strahlverwaltung, einer Art Schieberegister, in dem die Adressen der Parameter im Parameterspeicher verwaltet werden, wie in der Abbildung 7.6 mit 1,2,3 .. n angedeutet. Parallel zu den Adressen wird noch der Abstand zur Projektionsebene mit abgelegt.

Der Unterschied zu einem echten Schieberegister liegt in der Einsortierung. Und zwar werden die Strahlen nicht einfach eingeschoben, sondern nach dem Abstand zur Projektionsebene einsortiert. Strahlen, deren aktueller Abtastpunkt näher an der Projektionsebene liegt, werden dadurch früher von der Strahlverwaltung wieder ausgegeben.

Der Gesamt Ablauf über einen Strahl stellt sich folgendermaßen dar:

Die Adresse des aktuellen Abtastpunktes wird von der Strahlverwaltung an den Speichercontroller weitergegeben. Der Speicher liefert die Daten zur Interpolation und gleichzeitig die Sprungdistanz für Space-Leaping.

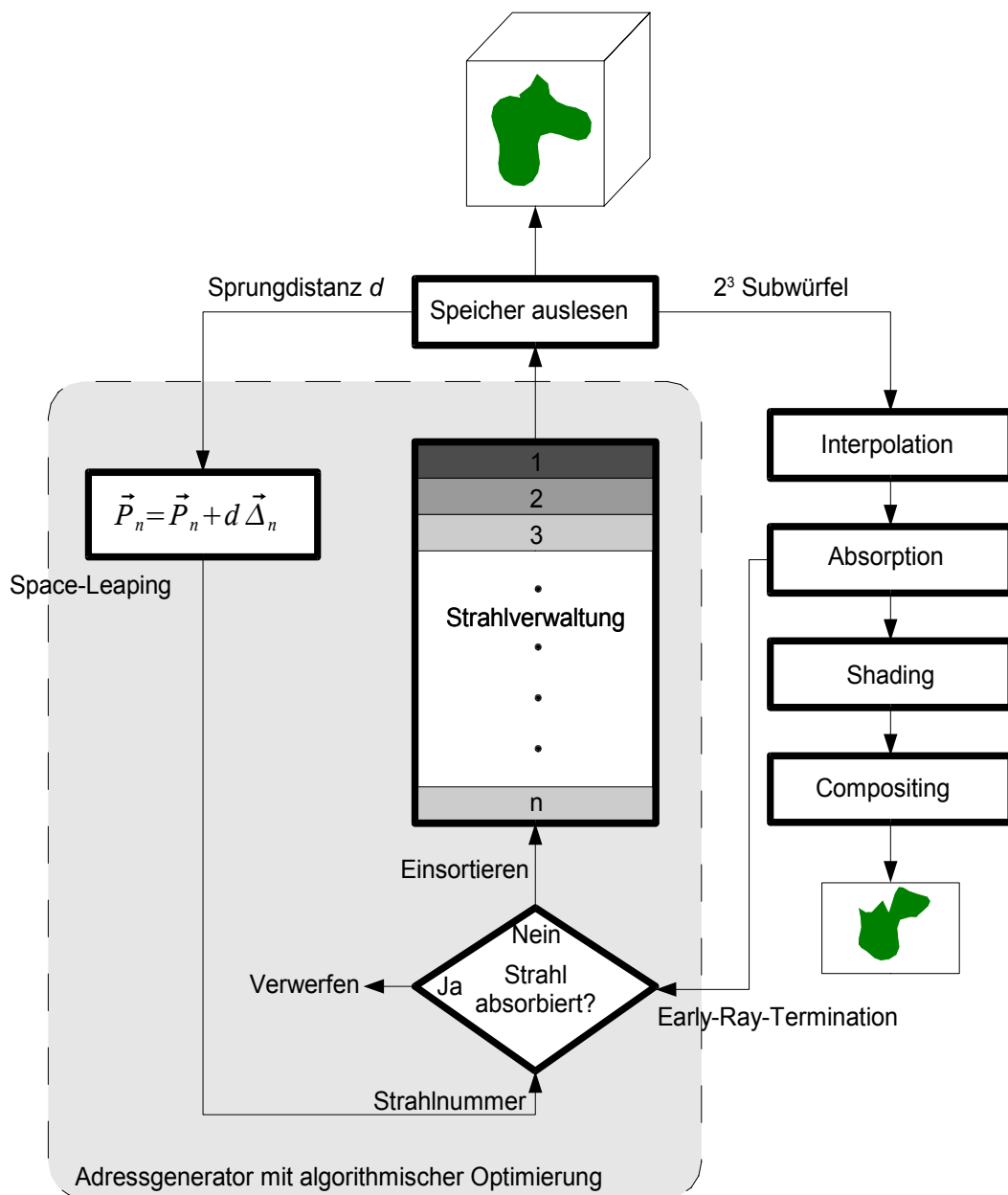


Abbildung 7.6: Pipeline mit erweitertem Adressgenerator für algorithmische Optimierung und Strahlverwaltung

Mit Hilfe der Sprungdistanz wird der nächste Abtastpunkt dieses Strahls berechnet und der neue Abstand zur Projektionsebene bestimmt. Parallel dazu wird der Abtastpunkt in die Ray-Casting-Pipeline geschoben.

Für Early-Ray-Termination muss eine Absorptionsberechnung zwischen Interpolation und Shading eingefügt werden. Hier wird bestimmt, ob die Intensität des Strahls unter einen Schwellwert gefallen ist.

Liegt das Ergebnis des nächsten Abtastpunktes für diesen Strahl vor, wird er von der Strahlverwaltung wieder einsortiert, falls er noch nicht absorbiert ist.

Da dieser ganze Vorgang mehrere Takte benötigt, wird in der Zwischenzeit mit jedem Takt ein weiterer Strahl von der Strahlverwaltung ausgegeben, bis alle Strahlen absorbiert sind, oder das Volumen verlassen haben.

Eine weitere Möglichkeit der Sortierung ist, nicht den Abstand zur Projektionsebene zu verwenden, sondern seine noch verbleibende Länge. Vom Host-System muss dann für jeden Strahl nicht nur der Anfangspunkt und der Distanzvektor bestimmt werden, sondern auch die maximale Länge des Strahls. Diese können so berechnet werden, dass der Strahl innerhalb des Volumens anfängt und nur bis zur Volumengrenze reicht, auch wenn die Projektionsebene außerhalb liegt.

Die Länge wird dann mit jedem Abtastpunkt oder Sprung reduziert. Die Sortierung ist dann umgekehrt, je größer der Wert, um so früher muss der Strahl wieder ausgegeben werden. Diese Variante hat den Vorteil, dass der Strahl nie außerhalb des Volumens liegt, wo auch keine Information für die Sprungdistanz verfügbar ist. Außerdem ist das Abbruchkriterium eines Strahls leichter zu überprüfen; es muss nur die verbleibende Länge auf Null getestet werden.

Nun stellt sich die Frage, wie aufwändig solch ein sortierendes Schieberegister als Strahlverwaltung zu realisieren ist.

7.4.2 Schaltungsrealisierung der Strahlverwaltung

Wie erwähnt, lässt sich die Strahlverwaltung durch ein abgewandeltes Schieberegister realisieren. Schieberegister haben die Aufgabe, ein Datenwort, das am Eingang angelegt wurde, mit jedem Takt um eine Stufe des Schieberegisters weiterzuschieben und so um die Anzahl der Stufen bis zum Ausgang zu verzögern. Eine Stufe eines Schieberegisters besteht nur aus einem Flip-Flop pro Bit des Datenwortes. Die Flip-Flops übernehmen alle ein Eingangs-Bit bei einer aktiven Taktflanke, falls ihr Freischaltsignal (Enable) *EN* aktiv ist.

Für unsere Anwendung besteht ein Datenwort aus einer Strahlnummer und aus einer Abstandsinformation zur Projektionsebene. Die Strahlnummer wird als Zeiger oder Adresse zu den eigentlichen Strahlparametern, wie aktuelle Position und Distanzvektor, in einem Parameter-RAM verwendet. Mit der Abstandsinformation wird verglichen, wo ein Strahl einsortiert wird. Je kleiner der Abstand, desto näher am Ausgang der Strahlverwaltung wird der Strahl einsortiert und somit bevorzugt bearbeitet.

Abbildung 7.7 zeigt drei Stufen der Strahlverwaltung. Den Ein- und Ausgängen einer Stufe des

Schieberegisters entsprechen die Datenbusse *data_in* und *data_out*. Ihnen wurde noch ein Gültigkeitssignal *shift_in* und *shift_out* hinzugefügt. Das Signal *shift_out* zeigt an, ob sich in der Stufe ein gültiger Strahl befindet und wird an *shift_in* der nächsten Stufe übergeben.

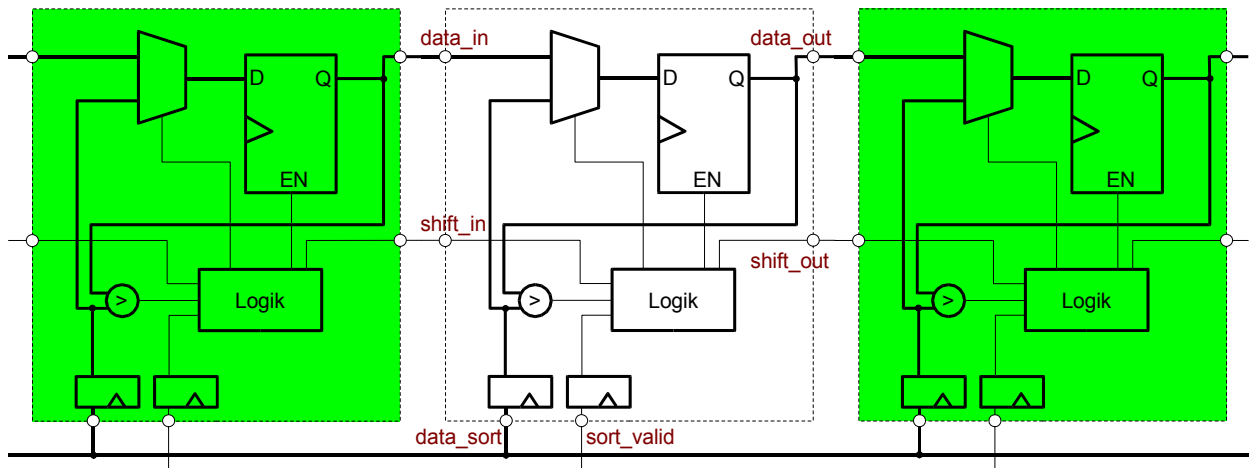


Abbildung 7.7: Schaltung eines kaskadierbaren Elements der Strahlverwaltung mit vorheriger und nachfolgender Stufe.

Weiterhin wurde das Schieberegister mit einem Bus als Sortiereingang *data_sort* mit seinem Gültigkeitssignal *sort_valid* erweitert. Beide müssen zur Entlastung des Sortier-Busses, der parallel an allen Stufen der Strahlverwaltung anliegt, in zusätzlichen Registern zwischengespeichert werden. Erst im nächsten Takt werden sie ausgewertet.

Ein Multiplexer kann, gesteuert durch eine Schaltlogik mit einem Komparator, wahlweise den Sortiereingang oder den Ausgang der vorherigen Stufe an den Register-Eingang *D* legen. Liegen keine gültigen Daten zum Sortieren an, wird die Stufe immer als Schieberegister geschaltet, also *data_in* übernommen.

Die Funktion der Strahlverwaltung besitzt folgenden Ablauf:

Zuerst wird die Strahlverwaltung mit gültigen, sortierten Strahlen gefüllt, indem an der Eingangsstufe die Daten nacheinander mit jedem Takt und aktivem *shift_in* angelegt werden. Der Sortiereingang ist über *sort_valid* inaktiv geschaltet. Wenn die Strahlverwaltung gefüllt ist, wird *shift_in* an der Eingangsstufe wieder abgeschaltet.

Mit einem in der Abbildung nicht gezeigten globalen Freischaltssignal (Enable) kann die Strahlverwaltung gestartet oder gestoppt werden, falls die Verarbeitung zum Beispiel auf den Speicher warten muss. Nach dem Starten arbeitet die Strahlverwaltung wie ein normales Schieberegister; an der Ausgangsstufe wird mit jedem Takt der nächste Strahl ausgegeben. Am

hinteren Ende der Strahlverwaltung werden über die *shift_in*- und *shift_out*-Signale die frei werdenden Stufen als ungültig markiert. Dies geschieht, indem *shift_in* gleich Null anzeigt, dass keine Daten vorhanden sind und die Null mit der nächsten Taktflanke an *shift_out* übernommen wird.

Die Strahlen werden nach der Ausgabe aus der Strahlverwaltung in der Ray-Casting-Pipeline weiterbearbeitet. Nach einigen Takten ist für den zuerst ausgegebenen Strahl der nächste Abtastpunkt und der neue Abstand zur Projektionsebene berechnet worden. Der neue Abstand mit der Strahlnummer wird dann über den Sortiereingang parallel an alle Stufen der Strahlverwaltung angelegt. Alle Stufen führen parallel einen Vergleich des neuen Abstandes mit dem in ihnen enthaltenen Abstandes zur Projektionsebene durch.

Alle Stufen, die feststellen, dass der Abstand ihres Strahles größer ist als der des einzusortierenden Strahls, schalten ihr *shift_out*-Signal ab und behalten ihren Wert für den nächsten Takt. Alle anderen Stufen, mit Abstand kleiner oder gleich, schieben die Daten normal mit aktivem *shift_out*-Signal weiter. Da die Daten sortiert in der Strahlverwaltung stehen, existiert nur eine Stufe, die den Abstand als kleiner gleich und das *shift_in*-Signal inaktiv erkennt. In dieser Stufe wird durch die Steuerlogik der Multiplexer so geschaltet, dass der einzusortierende Strahl übernommen wird.

Mit anderen Worten, die Strahlen mit größeren Abständen bleiben im Schieberegister stehen, während die anderen weitergeschoben werden. Der neu berechnete Strahl wird in die entstehende Lücke eingeschoben.

Um eine kontinuierliche Bearbeitung von Strahlen zu garantieren, muss die Strahlverwaltung mindestens so viele Stufen besitzen, wie man Takte benötigt um den nächsten Abtastpunkt zu berechnen. Erst wenn Strahlen beispielsweise durch Early-Ray-Termination nicht mehr einsortiert werden, kann die Strahlverwaltung leer laufen und es zu Lücken in der Bearbeitung kommen. Deshalb muss bei der Anzahl der Stufen ein Kompromiss gefunden werden zwischen Realisierbarkeit und möglichst vielen Stufen für die kontinuierliche Bearbeitung, selbst bei Reduzierung der Strahlen durch Early-Ray-Termination. Die optimale Anzahl der Stufen hängt allerdings, wie in Kapitel 9.1.2 erläutert wird, auch von der Organisation der verwendeten Speicherbausteine ab, da zu viele Strahlen wieder mehr Zeilenwechsel verursachen. Wie sich zeigen wird, ist eine Größe von 64 Strahlen, als 8x8 Strahlbündel organisiert, für den beschriebenen SDRAM-Baustein die beste Wahl.

7.5 Optimierung der Speicherarchitektur

Zur Speicherarchitektur gehört der physikalische Aufbau mit der Organisation der SDRAM-

Bausteine zu Speichermodulen und die Abbildung des dreidimensionalen Datenvolumens, das durch die Koordinaten X, Y und Z gekennzeichnet ist, auf den linearen Adressraum des Speichers.

Eine für die Bearbeitungsgeschwindigkeit optimale Speicherarchitektur für Ray-Casting, welche die aktiven Zeilen von SDRAM-Bausteinen als Zwischenspeicher ausnutzt, muss folgendes erfüllen:

1. Paralleles Auslesen aller Daten für die Interpolation:

Für eine gute Bildqualität sollten für die Interpolation des Abtastpunktes mindestens acht umgebende Volumenelemente, als 2^3 -Subwürfel, herangezogen werden. Um Mehrfachzugriffe zu vermeiden, müssen die acht Volumenelemente parallel ausgelesen werden können.

2. Richtungsunabhängigkeit:

Die Anzahl der Zeilenwechsel in den SDRAM-Bausteinen und die dadurch notwendigen Wartezeiten sollten von der Strahlrichtung weitgehend unabhängig sein. Dazu müssen die Daten würfelförmig in einer SDRAM-Zeile liegen.

3. Benachbarte Volumenelemente in unterschiedlichen Speicherbänken:

Bei der parallelen Abarbeitung mehrerer Strahlen (Multithreading) werden immer wieder lokal eng benachbarte Volumenelemente benötigt. Um bei jedem Strahlwechsel einen Zeilenwechsel zu vermeiden, dürfen benachbarte Volumenelemente sich nie in der gleichen SDRAM-Bank befinden.

Wären benachbarte Datenwürfel in der gleichen Bank, würde sogar ein Strahl alleine, der parallel zwischen den beiden Würfeln abtastet, bei jedem Abtastpunkt einen Zeilenwechsel verursachen, da für die Interpolation gleichzeitig Daten aus jedem Würfel benötigt werden.

Abbildung 7.8 zeigt eine mögliche Aufteilung der Daten, die diese drei Punkte erfüllt. Für die Interpolation müssen gleichzeitig acht Werte parallel zugreifbar sein, weshalb man acht getrennt adressierbare Speichermodule (A-H) benötigt. Jedes Volumenelement eines 2^3 -Subwürfels befindet sich in einem anderen Modul. Dies kann sehr einfach über die niederwertigsten Bit (LSB) der X-, Y- und Z-Koordinate jedes Volumenelementes dekodiert werden, wobei die Koordinaten in Binärdarstellung durch die Ziffernfolgen $z^n..z^0$, $y^n..y^0$ und $x^n..x^0$ mit dem Exponent 0 als niederwertigstes Bit repräsentiert werden (Tabelle 7.1):

z^0	y^0	x^0	Modul
0	0	0	A
0	0	1	B
0	1	0	C
0	1	1	D
1	0	0	E
1	0	1	F
1	1	0	G
1	1	1	H

Tabelle 7.1: Kodierung der Speichermodule mit den LSB der X, Y, Z-Koordinaten

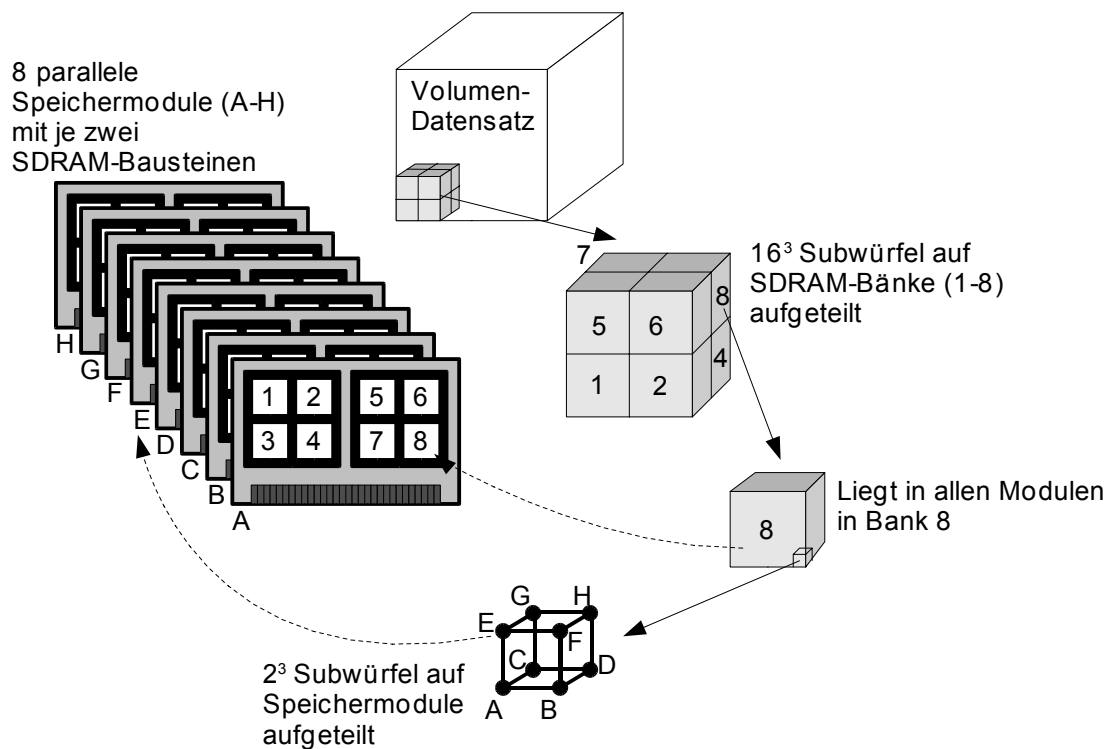


Abbildung 7.8: Aufteilung der Daten auf SDRAM-Speicher und Module

Für die Richtungsunabhängigkeit müssen auf die Zeilen der SDRAM-Bausteine möglichst große Datenwürfel verteilt werden. Kann ein SDRAM-Baustein beispielsweise 512 Datenworte in einer Zeile unterbringen, entspricht das einem 8^3 -Datenwürfel. Da die Daten aber noch auf 8 Module verteilt sind, entsprechen die Würfel, die in Abbildung 7.8 mit 1-8 gekennzeichnet sind, einem 16^3 Datenwürfel.

Die Abbildung eines Datenwürfels auf eine SDRAM-Zeile kann man erreichen, indem man die

Spaltenadresse der SDRAM-Bausteine aus den niedersten Bit der drei Koordinaten (X,Y,Z) zusammensetzt (x^0, y^0 und z^0 für die Moduladressierung ausgenommen). Beispielsweise besteht die Spaltenadresse einer Zeile mit 512 Datenworte aus 9 Adressbit. Zur Adressierung benötigt man von jeder Koordinate 3 Bit also (z^3, z^2, z^1) , (y^3, y^2, y^1) und (x^3, x^2, x^1) .

Damit nebeneinander liegende 16^3 Datenwürfel gleichzeitig in aktiven SDRAM-Bänken gehalten werden können, sind ebenfalls acht unterschiedliche Bänke pro Modul notwendig. Da SDRAM-Bausteine intern nur vier Bänke besitzen, müssen mindestens zwei SDRAM-Bausteine pro Modul vorgesehen werden. Die 16^3 Datenwürfel können, genau wie die 2^3 -Subwürfel, über drei Adressbit auf die zwei Bausteine mit je vier Bänken kodiert werden; zum Beispiel mit x^4 und y^4 zur Adressierung der internen Bänke und mit z^4 zur Auswahl einer der beiden Bausteine.

Alle höheren noch nicht verwendeten Bit der X,Y,Z-Koordinaten bilden die Zeilenadresse und sprechen gegebenenfalls auch weitere Speicherbausteine an. Zusätzliche SDRAM-Bausteine zur Speichererweiterung müssten gleichmäßig auf die Module verteilt werden.

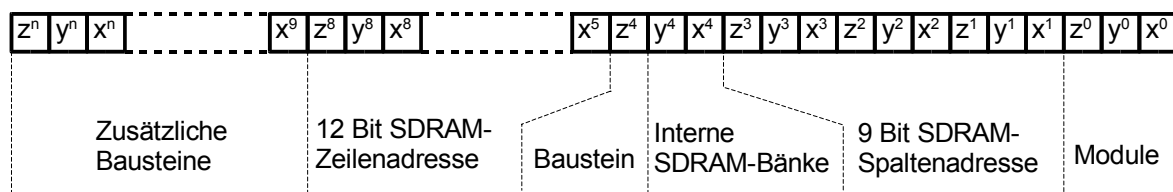


Abbildung 7.9: Gesamte Abbildung der X,Y,Z-Koordinaten auf den linearen Adressraum mit 128MByte SDRAM-Bausteinen

Abbildung 7.9 zeigt die komplette Zuweisung der X,Y,Z-Koordinaten auf den linearen Adressraum der Speichermodule und Bausteine. Innerhalb der SDRAM-Spalten- und Zeilenadresse ist ein Verschränken der drei Koordinaten nicht unbedingt notwendig, wird aber aus Symmetriegründen vorgeschlagen. Nur welche Bit aus den Koordinaten für die Spaltenadresse verwendet werden spielt eine Rolle, nicht die Reihenfolge.

7.6 Realisierungsmöglichkeit des Gesamtsystems

7.6.1 Grundsätzlicher Aufbau

Nachdem die grundlegenden Komponenten einer Volume-Rendering-Hardware mit algorithmischen Optimierungen und passender Speicherarchitektur beschrieben wurden, soll hier der Aufbau eines Gesamtsystems mit Anwendung und deren Zugriff auf die Rendering-Hardware beschrieben werden.

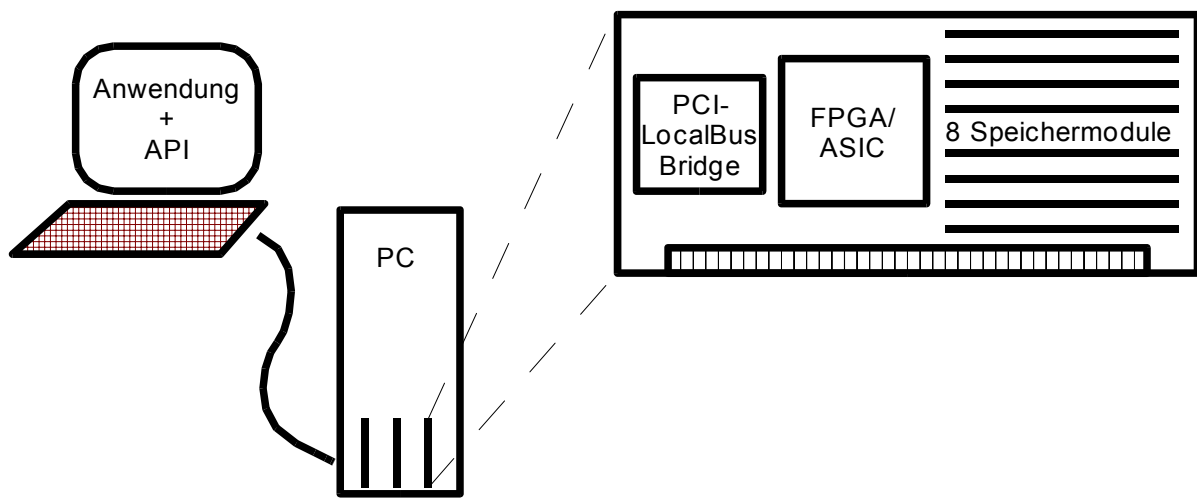


Abbildung 7.10: Realisierung der Rendering-Pipeline als PCI-Karte

Eine Anwendung die Volume-Rendering benötigt, muss die Hardware als Co-Prozessor verwenden und über bereitgestellte Software-Funktionen (API^I) auf die Rendering-Hardware zugreifen. Der Co-Prozessor übernimmt nur die Berechnung einzelner 8x8 Teilbilder. Die zugehörigen Startparameter und die Aufteilung des Gesamtbildes für jedes Teilbild muss die Anwendung mit Hilfe der API zur Verfügung stellen.

Sinnvoll ist ein Aufbau, der in Standard-PCs Verwendung finden kann, weshalb die Hardware als Steckkarte über einen PCI^{II}-Bus angeschlossen werden sollte (Abbildung 7.10). Aufgrund der großen Datenmengen, die übertragen werden müssen, kommt praktisch nur der PCI-Bus mit einer maximalen Übertragungsrate bei 133Mhz und 64Bit von 533MByte/s in Frage.

Abbildung 7.10 zeigt die Steckkarte, die 8 getrennte SDRAM-Speichermodule enthält. Hier können

^I API: Application Programming Interface

^{II} PCI: Peripheral Component Interconnect

auch für die Pre-Klassifikation notwendigen SRAM-Bausteine hinzugefügt werden. Die eigentliche Volume-Rendering-Berechnung mit den Speicherkontrollern finden in einem FPGA oder ASIC Platz. Um das Schaltungsdesign zu vereinfachen, kann die Kommunikation zwischen FPGA und PC über eine einfacher zu implementierende Local-Bus-Schnittstelle erfolgen. Ein zusätzliches Bauteil, bezeichnet als PCI/Local-Bus-Bridge, übernimmt dann die Umsetzung zum PCI-Bus. Möglich wäre auch eine direkte Implementierung der PCI-Schnittstelle, die von vielen Herstellern von FPGAs und ASICs als IP-Core^I angeboten werden.

Für die PCI/Local-Bus-Bridge spricht, dass bei FPGAs eine Programmierung über den PC möglich wird. Ansonsten muss die Programmierung über ein zusätzliches EPROM auf der Steckkarte realisiert werden.

7.6.2 Kommunikationsablauf

Im Folgenden soll erläutert werden, wie die komplette Berechnung eines Bildes im Zusammenspiel zwischen der Anwendung und dem Volume-Rendering-Co-Prozessor abläuft:

- Über die Anwendung wird zuerst ein Datensatz in den Volumenspeicher der Steckkarte übertragen.
- Über Benutzereingaben oder Standardwerte, werden die Parameter für Beobachterraichtung, Beleuchtung und Klassifikation eingestellt.
Diese Parameter bleiben für das gesamte Bild konstant und werden vor Beginn der Bildberechnung einmalig auf die Steckkarte übertragen.
- Bei jedem neuen Datensatz und immer wenn sich bei Klassifikationsänderungen die Leerräume im dargestellten Volumen ändern, muss die Distanzinformation neu berechnet werden. Dies kann in der Anwendung geschehen und muss dann in den Distanzspeicher geladen werden. Eine schnellere Variante ist die Integration der Distanzberechnung im FPGA (Kapitel 8.1).
- Die Anwendung muss die Projektionsebene in Kacheln der Größe 8x8 unterteilen.
- Für jeden der 64 Strahlen, die von der 8x8 Kachel der Projektionsebene ausgehen, wird der Start-, Offset- und Längenparameter berechnet. Hierbei muss ein *Clipping*[32] durchgeführt werden, d.h. Start- und Längenparameter müssen auf der Strahlgeraden so gewählt werden, dass sie innerhalb des Datenvolumens liegen. Sobald sich der Beobachter außerhalb des Objektes befindet, wird die Projektionsebene in den meisten Fällen nicht im

^I Intellectual Property, fertig platziert und verdrahtetes Bauteil, zum Einbau in ein ASIC oder Programmierung im FPGA.

Datenvolumen liegen. Gehen Strahlen am Datenvolumen vorbei, müssen sie nicht berechnet werden.

Der Adressgenerator erkennt zwar bei einem Startpunkt außerhalb, dass keine zugehörigen Daten vorhanden sind, kann aber dann nur schrittweise den Strahl weiter abtasten, bis Daten getroffen werden. Erst ab hier ist Distanzinformation enthalten und Space-Leaping kann leere Bereiche überspringen.

- Für jede Kachel werden die Parameter auf die Steckkarte geladen, die Berechnung gestartet und das Ergebnis des 8x8 Teilbildes ausgelesen.
- Durch doppelt angelegte Parameterspeicher kann die Visualisierung und das Herunterladen der Parameter parallel ablaufen, so dass die Zeiten für die Parameterberechnung und Kommunikation verdeckt werden können.
- Die Anwendung muss die Teilbilder zusammensetzen und darstellen.
- Erst nachdem neue Benutzereingaben gemacht wurden, geht der Zyklus von vorne los. Benutzereingaben können über eine Mauszeigerbewegung zum Beispiel eine Rotation des Objektes beschreiben, dann müssen mehrere aufeinander folgende Bilder erzeugt werden, bei denen sich nur die Beobachterraichtung ändert. Läuft dies für das menschliche Auge flüssig ab, ist das Ziel der Echtzeitdarstellung erreicht.

Da während diese Ablaufs, sehr viele Daten zwischen der Anwendung und der PCI-Karte ausgetauscht werden, müssen die Zeiten für die Übertragung geschätzt werden.

7.6.3 Abschätzung der Kommunikationszeiten

Zur Erzeugung eines kompletten Bildes müssen einige Parameter immer wieder für jedes 8x8 Teilbild auf die Steckkarte geladen werden. Tabelle 7.2 gibt eine Übersicht aller Parameter- und Datenspeicher der Hardware, mit Datenbreiten, Größen und den Übertragungszeiten, die erforderlich sind, um sie komplett beschreiben oder zu lesen. Die Übertragungszeiten sind für einen 133MHz, 64Bit PCI-Bus gerechnet, der eine Datenrate von 533MByte/s zulässt.

<i>Speicher</i>	<i>Datenbreite(Bit)</i>	<i>Kapazität (Worte)</i>	<i>Kapazität (Bytes gerundet)</i>	<i>Übertragungs- zeit (PCI 133)</i>	<i>Verwendung pro</i>
P(X,Y,Z)	3*23	64	576	1us	8x8 Teilbild
O(X,Y,Z)	3*15	64	384	0.7us	8x8 Teilbild
Ergebnis 8x8	8	64	64	0.1us	8x8 Teilbild
Strahllänge	10	64	128	0.2us	8x8 Teilbild
Reflectance-Map	32	2K	8k	15us	Bild
PR-LUT	8	64K	64K	120us	Bild
PO-LUT	8	256	512	1us	Bild
SDRAM-Distance	8	256k	256k	0.5ms (2.6ms)	Bild
SDRAM-Volume	16	1G	2G	4s (10s)	Anwendung

Tabelle 7.2: Parameter- und Datenspeicher

Die meisten Speicher werden nur zu Beginn der Visualisierung einmal für das gesamte Bild beschrieben und spielen deshalb für die Kommunikationszeiten keine große Rolle. Der Volumenspeicher wird normalerweise nur nach dem Starten der Anwendung geladen. Bei den SDRAM-Speichern sind deren Zugriffszeiten der begrenzende Faktor, weshalb die echten Ladezeiten in Klammer hinzugefügt wurden.

Weit mehr ins Gewicht fallen die 3 Parameter, die für jedes der 8x8 Teilbilder, als Initialisierung der Strahlen, übertragen werden müssen und das Ergebnisteilbild. Es handelt sich dabei um die Anfangspunkte P(X,Y,Z), den Richtungsvektor O(X,Y,Z) und die Strahllänge. Zusammengenommen muss etwa eine Größe von 1kByte für jedes Teilbild übertragen werden.

Tabelle 7.3 gibt für 3 Bildgrößen die summierten Übertragungszeiten aller zugehörigen 8x8 Teilbilder an. Die geschätzten Bildwiederholraten wurden für Volumenobjekte mit der gleichen Würfelkantenlänge, wie die Seitenlänge der Bildgröße ermittelt.

Man erkennt, dass die Kommunikationszeiten weit unter 10 Prozent der Bildberechnungsdauer liegt. Kann man diese Zeiten nicht tolerieren, können sie somit leicht durch ein doppeltes Auslegen der Parameterspeicher verdeckt werden. Während ein Parameterspeicher für die aktuelle Bildberechnung verwendet wird, kann die Anwendung die neuen Parameter berechnen, in den zweiten Speicher laden und das Ergebnisbild der vorherigen Berechnung auslesen. Bei der nächsten Berechnung vertauschen die Speicher ihre Rollen.

<i>Bildgröße</i>	<i>1024²</i>	<i>512²</i>	<i>256²</i>
<i>Anzahl der 8x8 Teilbilder</i>	16384	4096	1024
<i>Übertragungszeit</i>	35ms	8,9ms	2,2ms
<i>bei geschätzter Bildwiederholrate</i>	1Hz	6Hz	36Hz
<i>Übertragungszeit/Bild</i>	3,50%	5,30%	8,00%

Tabelle 7.3: Prozentualer Anteil der Übertragungszeit pro Gesamtbild mit PCI 133/64

Es sollte noch erwähnt werden, dass alle Berechnungen zwar mit der theoretischen Maximalgeschwindigkeit des PCI-Busses von 533MByte/s gerechnet wurden, diese aber in der Realität selten erreicht wird. Aber selbst mit der halben Geschwindigkeit, wären die Zeiten noch akzeptabel und könnten mit Speicherverdopplung verdeckt werden.

8 Ausbaufähigkeit der Architektur

8.1 Distance-Coding

Wie ein Kapitel 3.2 erläutert, berechnet der Distance-Coding-Algorithmus Sprungdistanzen zu den nächsten nicht leeren Volumenelementen und legt diese in einen eigenen Distanzspeicher. Der Distanzspeicher entspricht in der Organisation dem des Volumenspeichers, so dass mit der Adresse des Abtastpunktes der zugehörige Distanzwert ausgelesen werden kann. Man kann den Distanzspeicher auf ein Achtel reduzieren, indem die Distanzwerte für 2^3 -Subwürfel berechnet werden und nicht für jedes Volumenelement. Ein Subwürfel kann dann nur übersprungen werden,

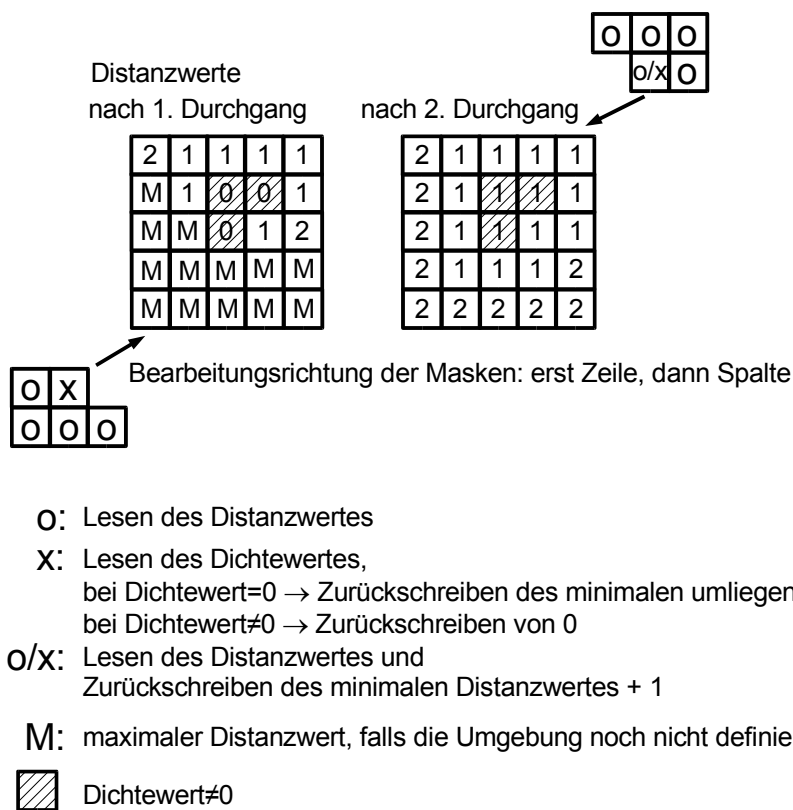


Abbildung 8.1: Ergebnis der Distance-Coding-Durchgänge

wenn er komplett leer ist.

Der Distance-Coding-Algorithmus kann in Software berechnet werden, so dass nur der fertige Datensatz in den Distanzspeicher geladen wird oder, wie im Anschluss erläutert, in Hardware umgesetzt werden.

Der Algorithmus kann durch Maskenoperationen mit zwei Durchgängen beschrieben werden. Die Maskenoperation vergleicht die räumlich um den aktuellen Berechnungspunkt liegenden

Distanzwerte und bestimmt daraus den Distanzwert. Abbildung 8.1 zeigt ein Beispiel innerhalb

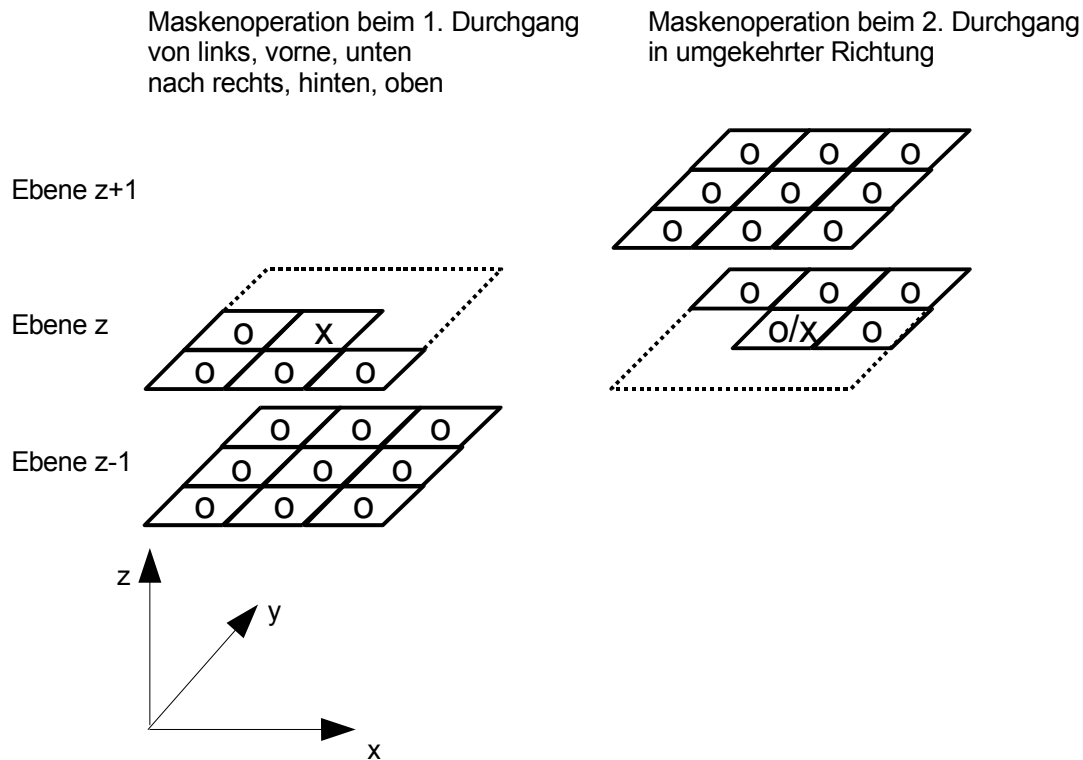


Abbildung 8.2: Bearbeitungsmasken von Distance-Coding mit zwei Durchgängen

einer Ebene mit den zugehörigen Masken. Dargestellt ist ein Ausschnitt aus dem Distanzspeicher, wobei an den schraffierten Kästchen im entsprechenden Teil des Volumenspeichers von Null verschiedene Dichtewerte stehen. Die Maske wird im ersten Durchgang mit der Schreibposition **X** zeilenweise von links nach rechts über die Speicherstellen geschoben. An jeder Schreibposition werden die umliegenden Distanzwerte gelesen und der minimale Distanzwert um eins erhöht in den Distanzspeicher geschrieben. Falls alle Distanzwerte undefiniert sind, wird der maximale Distanzwert **M** eingetragen. An den Positionen, an denen Dichtewerte ungleich Null vorliegen, wird als Distanzwert eine Null eingetragen.

Nach dem ersten Durchgang bleiben die Distanzen unter und links von nicht leeren Volumenelementen noch auf dem Maximalwert **M**. Erst nach dem zweiten Durchgang, der in umgekehrter Richtung abläuft, werden auch diese auf den richtigen Wert gesetzt. Im zweiten Durchgang können die Distanzwerte auf den Abtastabstand normiert werden und auch bei nicht leeren Volumenelementen auf einen Abtastabstand von eins gesetzt werden. Dadurch kann ohne Ausnahme der neue Abtastpunkt mit der Gleichung $\vec{P}_{n+1} = \vec{P}_n + \vec{\Delta} d$ berechnet werden, wie in Kapitel 3.2 erläutert.

Die Masken für den dreidimensionalen Fall sind in Abbildung 8.2 veranschaulicht. Beim ersten Durchgang müssen die neun Gitterpunkte eine Ebene tiefer unter der Schreibposition mit berücksichtigt werden, und beim zweiten Durchgang, die aus einer Ebene darüber.

Um den Algorithmus genau zu verstehen wird im Folgenden der Pseudo-Code gezeigt:

```
// 1. Durchgang
for z=0 .. zmax
  for y=0 .. ymax
    for x=0 .. xmax
      if Dichte=0
        minimum=MIN(
          dist[x-1,y,z],
          dist[x-1,y-1,z],dist[x,y-1,z],dist[x+1,y-1,z],
          dist[x-1,y-1,z-1],dist[x,y-1,z-1],dist[x+1,y-1,z-1],
          dist[x-1,y,z-1],dist[x,y,z-1],dist[x+1,y,z-1],
          dist[x-1,y+1,z-1],dist[x,y+1,z-1],dist[x+1,y+1,z-1])
        if minimum=undefined
          dist[x,y,z]=Max
        else
          dist[x,y,z]=minimum+1
      else
        dist[x,y,z]=0

// 2. Durchgang
for z=zmax..0
  for y=ymax..0
    for x=xmax..0
      if dist[x,y,z]=0 // Falls Volumenelement nicht leer
        dist[x,y,z]=1 // Abtastabstand von eins eintragen
      else
        minimum=MIN(
          dist[x+1,y,z],
          dist[x+1,y+1,z],dist[x,y+1,z],dist[x-1,y+1,z],
          dist[x+1,y+1,z+1],dist[x,y+1,z+1],dist[x-1,y+1,z+1],
          dist[x+1,y,z+1],dist[x,y,z+1],dist[x-1,y,z+1],
          dist[x+1,y-1,z+1],dist[x,y-1,z+1],dist[x-1,y-1,z+1])
        dist[x,y,z]=(minimum+1)/Abtastabstand //Normierung
```

Umsetzung in Hardware

Das Distance-Coding-Verfahren lässt sich relativ einfach hardwareunterstützt realisieren. Aus den zuvor beschriebenen Bearbeitungsmasken lassen sich folgende in Hardware einfach zu realisierende Funktionen ableiten:

- Überprüfen, ob ein Volumenelement leer ist
- Vergleichoperationen zwischen den benachbarten Distanzwerten
- Inkrementierung des kleinsten Distanzwertes um 1
- Zurückschreiben des neuen Wertes

Allerdings muss die Vergleichsoperation parallel 14 Distanzwerte zur Verfügung haben. Da die Bearbeitungsmaske aber in einer festgelegten Reihenfolge den Datensatz abarbeitet, können die Daten in ein Schieberegister geschoben werden. Wird das Datenvolumen in Säulen einer festen Grundfläche zerlegt, liegen die umliegenden Distanzwerte immer in der gleichen Verzögerungsstufe.

Abbildung 8.3 verdeutlicht die immer gleichbleibende Verzögerung der Daten in Bezug zu den Positionen der Bearbeitungsmaske, die grau hervorgehoben sind. Der Verzögerungswert 0

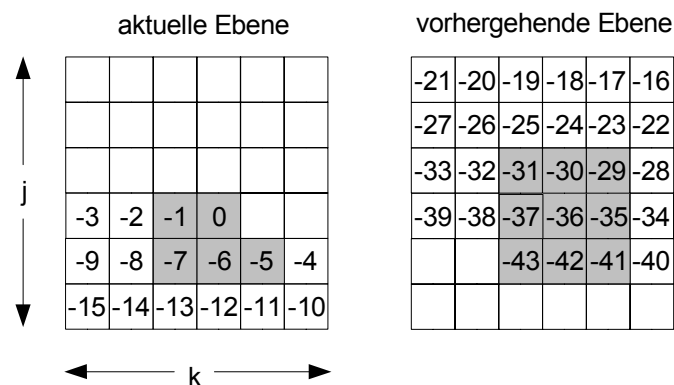


Abbildung 8.3: Bearbeitungsmaske und Verzögerung τ im Schieberegister bei fester Ebenengröße $j \times k$ (hier 6×6)

entspricht der aktuellen Schreibposition. Die vorherige Schreibposition befindet sich bei -1. Dessen Ergebnis kann aus der ersten Stufe des Schieberegisters ausgelesen werden und mit dem aktuellen Wert verglichen werden. Die weiteren Distanzwerte, mit denen in der aktuellen Ebene verglichen werden muss, befinden sich somit im 5., 6. und 7. Register. Ebenso wird mit den benachbarten Distanzwerten der vorherigen Ebene verfahren. Wird die Schreibposition eins weiter geschoben, rutschen auch die entsprechenden Daten im Schieberegister um eins weiter, so dass die zur

Verarbeitung notwendigen neuen Daten wieder an den gleichen Positionen im Schieberegister stehen.

Allgemein für eine Ebene mit k Spalten und j Zeilen befinden sich die Daten an folgenden Verzögerungsgliedern:

Verzögerungsglieder der aktuellen Ebene :

$$\tau^{-1} \quad \tau^0 \\ \tau^{-k-1} \quad \tau^{-k} \quad \tau^{-k+1}$$

und der vorhergehenden Ebene :

$$\tau^{-(j-1)k-1} \quad \tau^{-(j-1)k} \quad \tau^{-(j-1)k+1} \\ \tau^{-jk-1} \quad \tau^{-jk} \quad \tau^{-jk+1} \\ \tau^{-(j+1)k-1} \quad \tau^{-(j+1)k} \quad \tau^{-(j+1)k+1}$$

Es handelt sich bei diesen Verzögerungen nur um prinzipielle Werte. Bei der Realisierung muss berücksichtigt werden, dass für einen maximalen Durchsatz eine Vergleichsoperation nur 2 Werte pro Taktzyklus verarbeiten kann. Somit muss eine baumartige Struktur der Vergleiche aufgebaut werden. Die Zwischenergebnisse der Vergleiche werden über Register geführt, deren Verzögerung abgezogen werden muss.

Weiterhin lassen sich Vergleiche einsparen, indem die Vergleichsergebnisse der Dreiergruppen innerhalb einer Zeile verzögert werden, und nicht die Daten selbst. Dadurch kommt man mit 6 Vergleichen statt mit 14 aus; hierbei sind 2 für die Dreiergruppe und 4 für die Zeilen der aktuellen und vorhergehenden Ebene.

Ablauf über das Gesamtvolumen

Das eben beschriebene Verfahren verarbeitet nur einen Ausschnitt eines realen Datensatzes, da die Dimensionen der Ebenen festgelegt sind und nicht beliebig groß werden können. Wie viele dieser Ebenen gestapelt werden, ist unbegrenzt, so dass in der Höhe der Datensatz komplett durchlaufen werden kann. In den anderen Richtungen müssen die daraus entstehenden Säulen überlappend zusammen gesetzt werden, um den letzten Wert der vorhergehenden Säule zu berücksichtigen. Am Randbereich der Säulen darf kein Distanzwert geschrieben werden, da sich Teile der Masken außerhalb befinden und kein gültiger Wert berechnet wird. In Abbildung 8.4 sind unterschiedliche Säulen von oben gesehen, mit unterschiedlicher Schraffur und Grautöne für den Schreibbereich, dargestellt.

Wie beschrieben benötigt der Distance-Coding-Algorithmus zwei Durchgänge durch den Datensatz mit gegensätzlichen Richtungen. Es müssen erst alle Säulen in der ersten Richtung durchlaufen werden, bevor ein Richtungswechsel erfolgt. Hierbei müssen auch die Säulen nacheinander in der gleichen Bearbeitungsrichtung das Volumen durchlaufen, wie es auch innerhalb der Ebene mit den

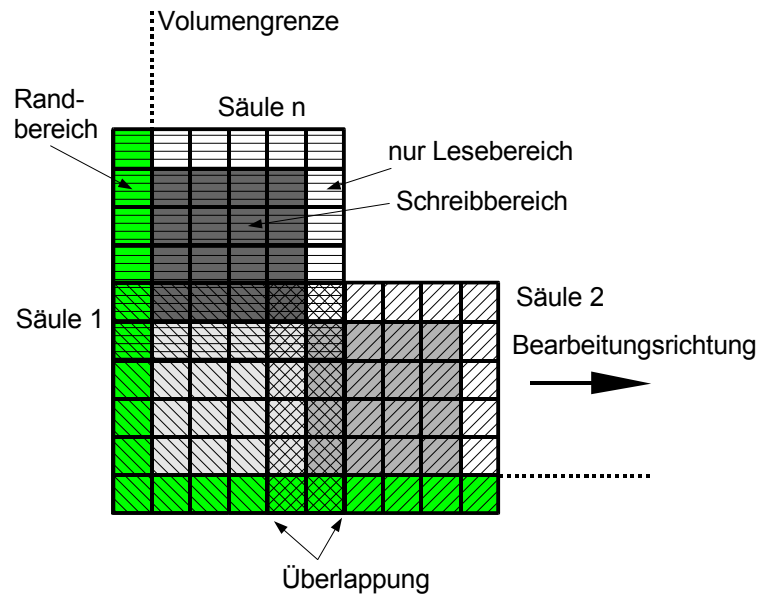


Abbildung 8.4: Bearbeitungssäulen von oben gesehen mit Überlappung und Randbereich

Masken geschieht. Also beispielsweise im ersten Durchgang von links nach rechts die Säulen aneinander reihen, dann von vorne nach hinten die nächsten Säulenreihen und innerhalb der Säule von unten nach oben. Beim 2. Durchgang ist die Bearbeitungsrichtung genau entgegengesetzt. Die Bearbeitungsrichtung wird nur durch die Reihenfolge festgelegt, mit welcher der Adressgenerator die Volumenelemente auswählt und in das Schieberegister gibt.

Das Schaltwerk mit den Vergleichern muss für den 2. Durchgang nicht verändert werden. Es reicht die Daten in umgekehrter Richtung einzuschieben. Vor Beginn jedes Durchgangs müssen alle Registerstufen auf den maximalen Distanzwert gesetzt werden, damit keine Vergleiche mit alten, ungültigen Werten zu falschen Ergebnissen führt.

Als Spezialfall muss an den Rändern des Gesamtvolumens für außerhalb liegende Maskenwerte ebenfalls der maximale Distanzwert angenommen werden, wie in Abbildung 8.4 grün unterlegt angedeutet.

Die Funktion des gesamten Verfahrens wurde durch ein synthetisierbares VHDL-Programm überprüft. Da dieser Block unabhängig vom Ray-Casting-Algorithmus verwendet werden kann, gibt es zwei Möglichkeiten der Realisierung im FPGA. Falls im FPGA noch genügend Platz vorhanden ist, kann das Distance-Coding parallel mit implementiert werden. Nur bei Bedarf wird es dann durch ein Steuerregister aktiviert.

Eine weitere Möglichkeit ist die Rekonfiguration des FPGAs für Distance-Coding. Es kostet zwar einige Sekunden, den FPGA neu zu programmieren, dafür steht aber bei einem stark belegten FPGA

mehr Platz für den Ray-Casting-Algorithmus zur Verfügung. Der Distance-Coding-Block benötigt neben einer kleinen Steuerung nur noch den Speicherkontroller.

8.2 Perspektive

Mit der beschriebenen Rendering-Architektur lässt sich die perspektivische Darstellung für kleinere Sehwinkel ohne Änderung realisieren. Da die Abtastpunkte entlang des Strahls durch Aufsummieren eines Differenzvektors berechnet werden, muss nur für jeden Strahl, der Vektor neu

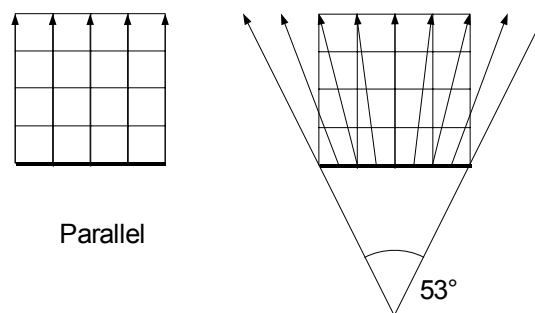


Abbildung 8.5: Vergleich Parallel- und 53°-Projektion

bestimmt werden. Es muss allerdings berücksichtigt werden, dass die Strahlen auseinander gehen. Soll keine Information verloren gehen, muss der Abtastabstand in Höhe und Breite der Projektionsebene verringert werden. In der Tiefe kann der Abtastabstand von 1 beibehalten werden. Abbildung 8.5 zeigt bei einem würfelförmigen Objekt, dass bei einem Perspektivenwinkel von 53° und halbiertem Abtastabstand, an der Rückseite des Würfels noch genau ein Abstand von 1 der Strahlen erreicht wird. Wie die Simulation zeigt, wird die Rechenzeit dadurch nur um etwa 10-20% erhöht, da ein großer Teil der Strahlen den Würfel bereits zu Anfang verlässt.

Allgemeiner wird der Zusammenhang zwischen Perspektivenwinkel und Abtastabstand in der folgenden Abbildung 8.6 gezeigt. Es wird untersucht, mit welchem Abtastabstand in der Projektionsebene begonnen werden muss, damit bei der Durchquerung eines Würfels an der Rückseite noch ein Abtastabstand von 1 eingehalten wird. Hier kann man wieder bei etwa 53° den Abtastabstand von 0.5 ablesen.

Über den Strahlensatz lässt sich der notwendige Abtastabstand in der Projektionsebene ($s_{x,y}$) in Abhängigkeit vom Perspektivenwinkel β herleiten zu:

$$s_{x,y} = \frac{1}{2 \tan\left(\frac{\beta}{2}\right) + 1}$$

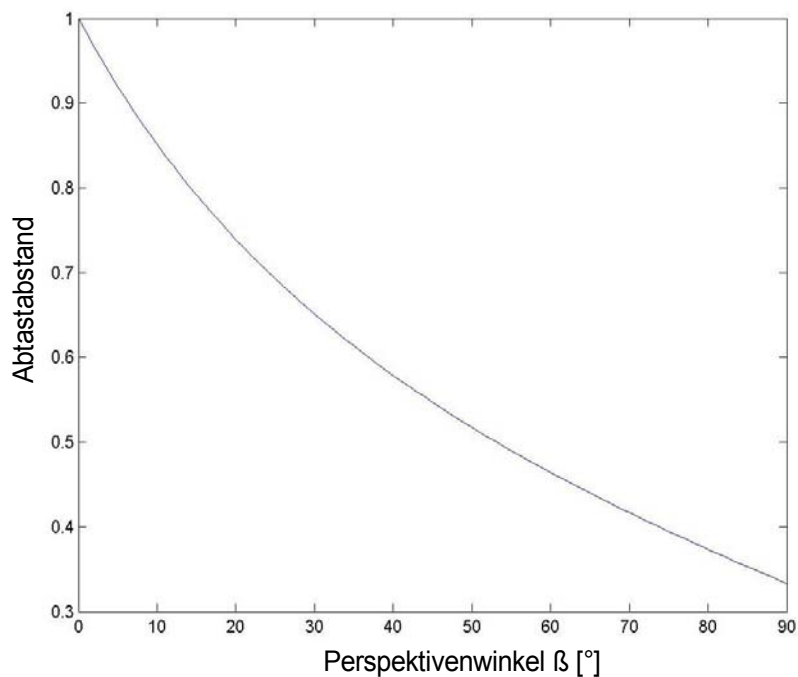


Abbildung 8.6: Notwendiger Abtastabstand $s_{x,y}$ in der Projektionsebene zur Durchquerung eines Würfels, in Abhängigkeit vom Perspektivenwinkel β .

Das Diagramm in Abbildung 8.7 zeigt, wie sich bei gegebenem Abtastabstand ($s_{x,y}$) von 0.5 die Eindringtiefe des Sehstrahls im Verhältnis zur Breite der Projektionsebene ändert. Bei 53° erkennt man wieder das Verhältnis von 1, was einem Würfel entspricht. Wird der Perspektivenwinkel allerdings auf 90° vergrößert, dürfen die Sehstrahlen nur noch halb so tief eindringen.

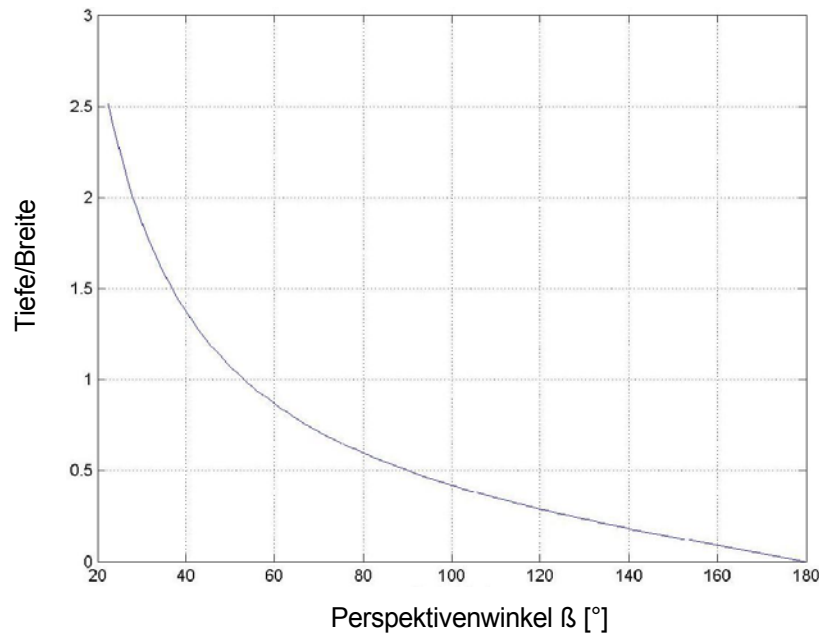


Abbildung 8.7: Verhältnis der Eindringtiefe des Sehstrahls zur Breite der Projektionsebene bei einer Abtastabstand von 0.5

Das Diagramm wurde berechnet mit:

$$\frac{\text{Tiefe}}{\text{Breite}} = \frac{1 - s_{x,y}}{2 s_{x,y} \tan\left(\frac{\beta}{2}\right)}$$

Bei größeren Winkeln oder einem Tiefen zu Breiten Verhältnis größer 1, müssen daher eventuell zusätzliche Projektionsebenen eingefügt werden. Die Anwendung muss die Strahlen so initialisieren, dass sie von einer zur nächsten Projektionsebene laufen und die entstehenden Zwischenbilder miteinander durch Compositing^I verrechnen. Hierbei erhöht sich der Abtastabstand in Breite und Höhe und es ändert sich der Perspektivenwinkel für die Differenzvektorberechnung von einer zur nächsten Ebene, wie in Abbildung 8.8 zu sehen.

^I vergleiche hierzu Abschnitt Compositing in Kapitel 8.3

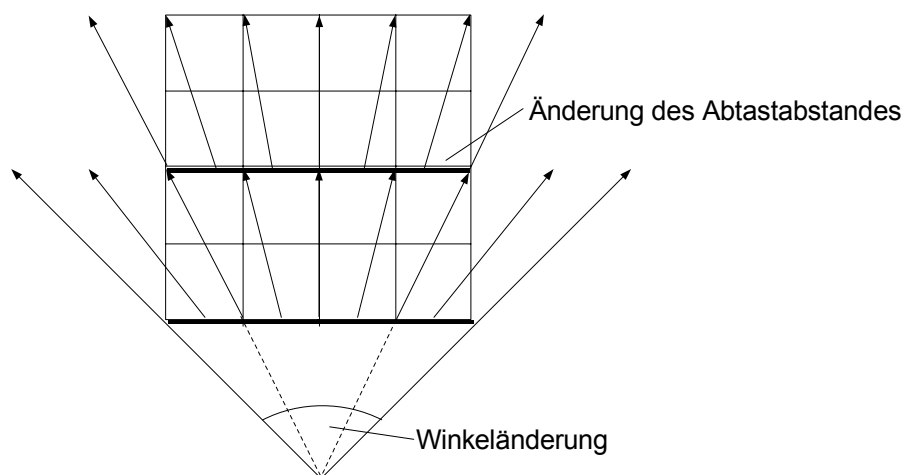


Abbildung 8.8: Einfügen von Zwischenebenen bei großen Perspektivenwinkeln

Mit anderen Worten ausgedrückt, werden, sobald sich die Strahlen um mehr als 1 voneinander entfernen, neue Strahlen eingefügt.

Die Anwendung kann durch Auswertung des Registers für terminierte Strahlen, das durch Early-Ray-Termination gesetzt wird, die Zahl der Strahlen von Ebene zu Ebene weiter reduzieren.

8.3 Erweiterungsmöglichkeiten durch die Anwendungs-Software

Multi-Volume

Aus Sicht der Anwendungs-Software wird die Ray-Casting-Hardware als Funktionsaufruf eingebunden, der ein 8x8-großes Teilbild zurückgibt. Die Anwendung muss neben einigen anderen Parametern vor allem den Beobachtervektor zur Verfügung stellen. Durch ihn ist die Rotation und Translation des zu visualisierenden Objektes bestimmt.

Eine Erweiterungsmöglichkeit, die durch die Anwendung realisiert werden kann, ist das Verwalten mehrerer Volumenobjekte in einer Szene. Dies wird oft durch den Begriff *Multi-Volume* bezeichnet (siehe Abbildung 8.9).

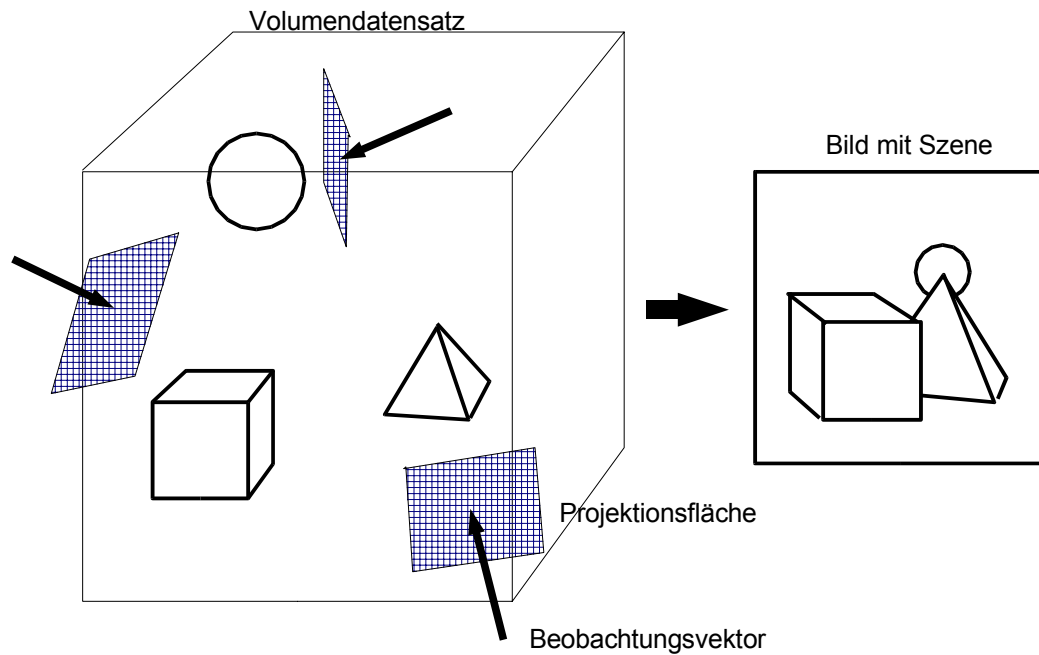


Abbildung 8.9: Zusammensetzung einer Szene aus unabhängigen Volumenobjekten

Die Anwendung muss die Position der Objekte zueinander festlegen. Die Objekte können parallel im Volumenspeicher liegen, ohne Rücksicht auf ihre Lage in der zu visualisierenden Szene. Die Objekte werden getrennt aus unterschiedlichen Richtungen visualisiert und die dabei entstehenden Einzelbilder zur Szene zusammengesetzt. Wird ein Objekt bewegt, muss nur dessen Beobachtungsvektor verändert werden und das neu entstandene Bild mit den schon erzeugten Bildern der nicht bewegten Objekte zusammengesetzt werden.

Für das Zusammensetzen überlappender Bilder gibt es mehrere Möglichkeiten, die in Oberflächengraphik ebenfalls unter dem Namen *Compositing* bekannt sind [32]. Die am meisten verwendete entspricht auch dem Compositing entlang des Sehstrahls in der Volumenvisualisierung und wird durch den **over**-Operator beschrieben. Wenn Bild A vor Bild B liegt, wird das Ergebnisbild in jedem Pixel berechnet durch:

$$A \text{ over } B = C_A \alpha_A + C_B (1 - \alpha_A)$$

Hierbei stehen C_A und C_B für die Farbe des jeweiligen Pixels und α_A für die Transparenz des Pixels in Bild A. Durch diese Berechnung werden auch halbtransparente Objekte berücksichtigt, wodurch auch hinten liegende Objekte durchscheinen können.

Multi-Board

Auf die gleiche Weise kann eine Anwendung auch mehrere Platinen oder Volume-Rendering-Einheiten verwalten. Dies wird als *Multi-Board* bezeichnet. Um die volle Leistung zu erhalten,

muss nur sichergestellt sein, dass im PC-System der PCI-Bus die Datenrate liefern kann (vergleiche hierzu Kapitel 7.6). Eventuell muss ein System mit mehreren PCI-Bussen verwendet werden.

Das Gesamtsystem hat dadurch mehrere getrennt ansprechbare Volumenspeicher zur Verfügung. Es können unterschiedliche Objekte, aber auch ein einzelnes zu großes Objekt, auf die Volumenspeicher aufgeteilt und parallel berechnet werden. Auf jeder Rendering-Einheit gibt es eigene Projektionsebenen, die wie zuvor beschrieben zu einer Szene zusammengesetzt werden können.

8.4 Farbe und zusätzliche Klassifizierung

Bei den von uns betrachteten Datensätzen handelt es sich um dreidimensionale Skalarfelder, die somit nur einen Wert, wie beispielsweise die Dichte, einem Punkt im Raum zuordnen. Dadurch haben die darzustellenden Objekte von sich aus keine Farbe. Zur besseren Unterscheidung der Dichtebereiche können über die Klassifikation dennoch Farbwerte^I zusätzlich zur Opazität (Alpha-Wert α) zugeordnet werden. Genauso können zusätzlich die Koeffizienten für den ambienten, diffusen und spiegelnden Lichtanteil zugeordnet werden, um unterschiedliche Materialeigenschaften zu simulieren. Um den Aufwand für die Interpolation der Zusatzangaben zu umgehen, werden diese normalerweise in der Post-Klassifikation eingefügt und müssen dann nur im Shading und Compositing berücksichtigt werden.

Eine Möglichkeit der Implementierung ist in [44] beschrieben (Abbildung 8.10). Hier wird jedem Abtastpunkt mit der Post-Klassifikation über die Dichte die Werte $R, G, B, \alpha, k_a, k_d$ und k_s zugeordnet. Der ambiente, diffuse und spiegelnde Lichtanteil wird wie in Kapitel 5.3 beschrieben in der Illumination berechnet. Deren Realisierung ist ähnlich der in Kapitel 6.7 gezeigten, nur dass die drei Reflexionsbestandteile getrennt in der Reflexionstabelle abgelegt sind, damit sie mit den Koeffizienten multipliziert werden können.

Im Combiner wird bei dieser Implementierung nur der diffuse Lichtanteil mit den drei Farben R, G, B verrechnet:

I Farbwerte: RGB für Rot, Grün, Blau

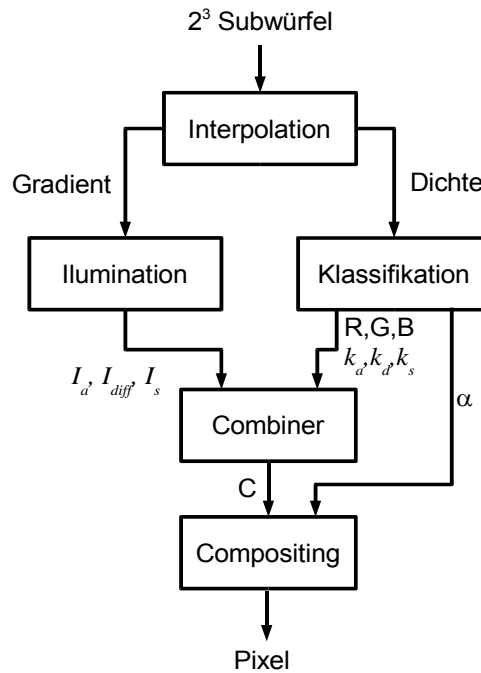


Abbildung 8.10: Farbberechnung bei VIZARD II

$$C_{Farbe} = k_a I_a + k_d I_{diff} \cdot Farbe + k_s I_s$$

Im Compositing wird die Pixel-Farbe (R,G,B,α) für jeden Kanal getrennt, wie in Kapitel 5.5 für einen Kanal beschrieben errechnet, wobei die Farbkanäle R,G,B auch mit α gewichtet und mit 1-α abgeschwächt werden:

$$C_{n+1} = C_n + C_{Sample} \alpha_{Sample} (1 - \alpha_n)$$

$$\alpha_{n+1} = \alpha_n + \alpha_{Sample} (1 - \alpha_n)$$

C_{Sample} und α_{Sample} entsprechen der Farbe und der Opazität des aktuellen Abtastpunktes.

8.5 Einbeziehen von Polygondaten

Die meisten Objekte der realen Welt lassen sich in der Computergrafik wesentlich leichter, realistischer und ressourcenschonender mit Polygondaten darstellen. Auch ist von vielen Objekten das Innenleben nicht bekannt, oder für den Anwendungsfall uninteressant, weshalb der Wunsch besteht, im Normalfall die Polygondarstellung zu wählen und nur wo es notwendig und sinnvoll ist auf Volumendaten zurückzugreifen. Ein Volume-Rendering-System alleine, ohne die Möglichkeit Oberflächenobjekte in Polygondarstellung mit einzubinden, wird wenig Akzeptanz finden. Eine mögliche Anwendung ist zum Beispiel die Computer-unterstützte Operationsplanung, bei der das zu operierende Organ nach einer Computertomographie oder Magnetresonanz-Aufnahme als

Volumenobjekt vorliegt und der Arzt mit virtuellen medizinischen Werkzeugen die Operation am Bildschirm proben kann. Die Werkzeuge lassen sich wesentlich realistischer als Oberflächengrafik darstellen.

Ein Problem, das bei der Kombination von Oberflächen- und Volumengrafik zu lösen ist, ist das Verdeckungsproblem. Wenn der Arzt beispielsweise mit seinem virtuellen Skalpell in das Gewebe eintaucht, muss durch Verschwinden der Schneide erkennbar sein, wann das Skalpell die richtige Schnitttiefe besitzt. Das Volumenobjekt muss das Oberflächenobjekt verdecken können, aber genauso auch umgekehrt, muss der Griff des Skalpells undurchsichtig sein und vor dem Volumenobjekt liegen.

Um das Problem aus Sicht eines Volume-Rendering-Systems zu lösen, benötigt man entlang eines Sehstrahls die Information, in welcher Entfernung ein Oberflächenobjekt kreuzt und muss dieses im Compositing berücksichtigen. Diese Entfernung entspricht dem Abstand zur Projektionsebene. Bei unserem Multithreading-Ray-Casting-System finden wir diese Information zur Strahlsortierung bereits im Adressgenerator (Kapitel 7.4). Die Abstandsinformation kann über Schieberegister zum richtigen Zeitpunkt mit dem verarbeiteten Abtastpunkt und dessen Reflexion an die Compositing-Stufe gegeben werden. Im Compositing ist auch bekannt, für welchen Bildpunkt ein Beitrag berechnet wird.

Wenn man sich die OpenGL-Pipeline aus Kapitel 2.2 anschaut, verfügt man über die gleiche Information, nämlich Abstandswert im Z-Puffer und Farbe mit Helligkeit jedes Bildpunktes, nach den Fragmentoperationen im Frame-Puffer. Wenn man diese Daten ebenfalls an die Compositing-Stufe der Ray-Casting-Pipeline gibt, kann dort für jeden Sehstrahl mit zugehörigem Bildpunkt, durch eine einfache Vergleichsoperation, entschieden werden, welche Information zuerst berücksichtigt werden muss. Den prinzipiellen Aufbau zeigt Abbildung 8.11.

Wird ein Bildpunkt aus dem Frame-Puffer ausgewählt, kann dessen Farbe mit der Compositing-Gleichung aus Kapitel 8.4 eingerechnet werden. Ist der Bildpunkt undurchsichtig, kann der zugehörige Sehstrahl sofort beendet und Early-Ray-Termination ausgelöst werden.

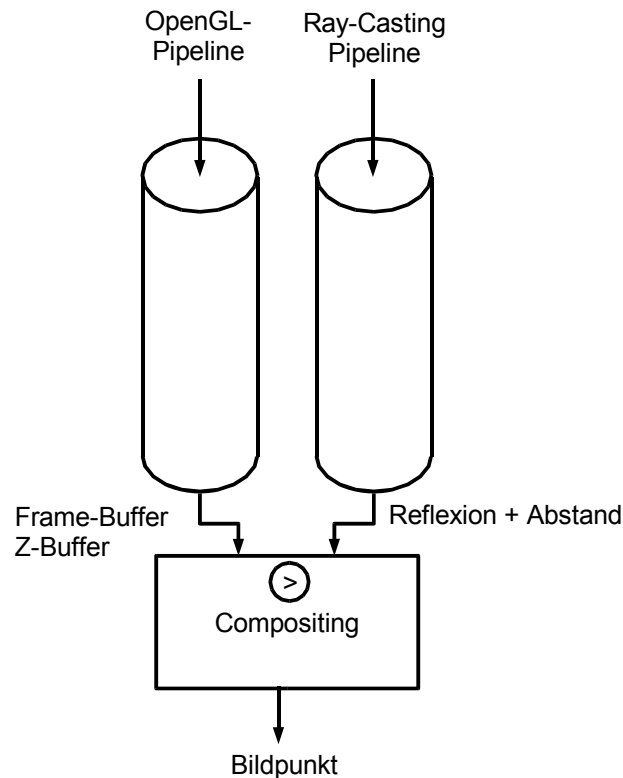


Abbildung 8.11: Abstandsvergleich zwischen Volumen und Oberflächendaten im Compositing

8.6 Schatten

Wie in Kapitel 2.3.5 kurz erläutert entstehen Schatten beim Volume-Ray-Tracing-Verfahren durch Berücksichtigung einer seitlichen Lichtquelle mit Absorption. Kajiya und von Herzen [54] zeigten, dass dies mit einem zusätzlichen Durchlauf für jede seitliche Lichtquelle möglich ist.

Für eine seitliche Lichtquelle wird im ersten Durchlauf die Intensität I_i entlang der Lichtstrahlen an jedem Volumenelement bestimmt und in einem zugehörigen Puffer $I_s(X,Y,Z)$ zwischengespeichert. Die Intensität wird näherungsweise mit der Absorptionsgleichung entlang des Lichtstrahls bestimmt:

$$I_{i+1} = I_i \cdot (1 - O_i)$$

Da eine Absorption berechnet wird und leere Volumenelemente keine Reflexion verursachen, kann auch bei der Berechnung der seitlichen Lichtquelle Early-Ray-Termination und Space-Leaping berücksichtigt werden.

Im zweiten Durchlauf wird das normale Ray-Casting durchgeführt. Der einzige Unterschied liegt im Compositing, da hier die Intensität der seitlichen Lichtquelle an jedem räumlichen Punkt aus dem Schattenpuffer $I_s(X,Y,Z)$ zusätzlich berücksichtigt werden muss:

$$I_{out,i} = I_{phong,i} O_i I_s(X, Y, Z) + I_{out,i-1} (1 - O_{i-1})$$

Dieses Verfahren setzt allerdings einen Schattenpuffer voraus, der die gleiche Größe besitzt, wie der Volumenspeicher. Dies wäre realisierbar, indem man die Datenbusbreite des Volumenspeichers erhöht, und die Intensität der seitlichen Lichtquelle in den Volumenspeicher bei dem zugehörigen Volumenelement zurückschreibt. Allerdings wird dadurch nur ein Drittel der Bildwiederholrate des reinen Ray-Casting erreicht, da hier dreimal auf den Speicher zugegriffen werden muss. Das erste mal für das Lesen der Opazitäten, zweitens für das Zurückschreiben der Intensität in den Schattenpuffer und drittens das nochmalige Lesen für das Ray-Casting-Durchgang.

Ohne Erweiterung des Volumenspeichers kommt man aus, wenn die Berechnung der seitlichen Lichtquelle innerhalb einer Strahlbündelberechnung vorgenommen wird. Dadurch ist ein kleinerer Schattenpuffer für das gesamte Strahlbündel notwendig. Die Größe liegt bei 64kByte^I, wenn mit 64 Strahlen, 1024 als maximale Länge und ein Byte für die Intensität gerechnet wird.

Werden beliebige Richtungen der seitlichen Lichtquelle zugelassen, muss über Zwischenspeicher die noch bestehende Intensität der seitlichen Lichtquelle von einem Strahlbündel zum nächsten weitergegeben werden. Notwendig ist einmal ein Zwischenspeicher zum Nachbarstrahlbündel mit der Größe 8x1024 und ein Zwischenspeicher zu der nächsten Ebene (1024x1024) von Strahlbündeln. Die Zwischenspeicherung der Ebene kann bei einer FPGA-Implementierung allerdings nur extern erfolgen.

Auf diese Zwischenspeicher kann man verzichten, wenn man ähnlich wie bei VIRIM (Kapitel 4.3) nur seitliches Licht innerhalb der Ebene zulässt, so dass nur Abhängigkeiten zum direkt daneben liegenden Strahlbündel bestehen. Wenn die Strahlbündel in der gleichen Richtung abgearbeitet werden, wie das seitliche Licht einfällt, kann sogar auf den Zwischenspeicher zwischen den Strahlbündeln verzichtet werden, wie in Abbildung 8.12 angedeutet. Hier im Beispiel kommen die seitlichen Lichtstrahlen von links vorne. Nach der Berechnung der seitlichen Lichtstrahlen bis an den rechten Rand des Strahlbündels, stehen dort im Schattenpuffer die Intensitäten, wie sie für das nächste Strahlbündel links zum weiterrechnen benötigt werden. Man muss somit nur den Schattenpuffer in einer Richtung abwechselnd für jedes Strahlbündel spiegeln.

I Ein maximales Volumen von 1024³ angenommen

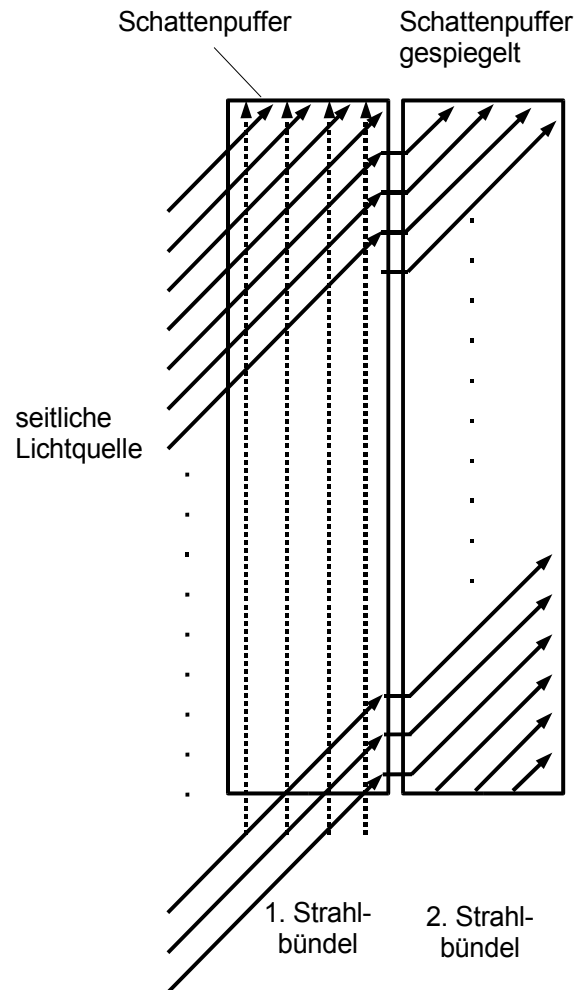


Abbildung 8.12: Ray-Tracing mit gespiegelmtem Schattenpuffer

Diese Spiegelung kann durch eine einfache Umrechnung einer Koordinate realisiert werden. Der Schattenpuffer wird mit den relativen Koordinaten X, Y, Z im Strahlbündel adressiert. Da er nur die Dimensionen $8 \times 8 \times 1024$ hat, sind für die X - und Y -Koordinaten nur drei Bit notwendig. Eine Spiegelung in X -Richtung lässt sich durch eine einfache Negation der drei Bit realisieren:

<i>normal</i>	<i>gespiegelt</i>	<i>normal Binär</i>	<i>gespiegelt Binär</i>
0	7	000	111
1	6	001	110
2	5	010	101
..
7	0	111	000

Der Adressgenerator muss so ausgelegt werden, dass innerhalb eines Strahlbündels zuerst die

Absorption des seitlichen Lichts und dann das normale Ray-Casting berechnet wird. Trotz der Bearbeitung innerhalb eines Strahlbündels werden hier alle Daten zwei mal ausgelesen, wodurch sich die Bildwiederholrate gegenüber dem reinen Ray-Casting halbiert. Die Bildwiederholrate könnte noch etwas erhöht werden, wenn die Abarbeitung noch die Datenkohärenz berücksichtigt. Gehen beide Durchläufe nacheinander komplett durch das Strahlbündel, müssen alle benötigten Subwürfel zwei mal geladen werden. Dies kann vermieden werden, wenn beide Berechnungen innerhalb eines gerade geladenen Subwürfels durchgeführt werden. Ob sich der Aufwand hierfür lohnt und wie sich dies auf die Komplexität der Steuerung und des Adressgenerators auswirkt, müsste noch untersucht werden.

Das VIRIM-System kann das seitliche Licht und den Ray-Casting-Schritt in einem Durchgang berechnen. Dies wird durch die Umsetzung des Heidelberger Ray-Tracing-Algorithmus erreicht [85]. Man beschränkt sich hier auf eine 45°- und eine 0°-Lichtquelle. Die Strahlen des 45°-Lichtes werden so bestimmt, dass ihre Abtastpunkte mit denen des 0°-Lichtes zusammenfallen. Durch die kontinuierliche Abarbeitung der Abtastpunkte von links nach rechts und parallel zur Projektionsebene, kann daher für jeden Abtastpunkt gleichzeitig die Absorption des 45°-Lichtes und das Compositing mit dem 0°-Licht berechnet werden.

Diese kontinuierliche Abarbeitung der Strahlen ist aber bei der hier vorgestellten Multithreading-Architektur mit Sortierung der Strahlen nicht gegeben, weshalb sich das Heidelberger-Ray-Tracing-Verfahren nicht direkt umsetzen lässt. Speziell für das Ray-Tracing-Verfahren könnte man allerdings auf die Sortierung verzichten und nimmt eine erhöhte Anzahl von Zeilenwechsel in Kauf.

8.7 Deformation

Um eine virtuelle Operation möglichst realistisch darstellen zu können, ist es notwendig, dass sich beispielsweise Gewebe beim Auftreffen eines Skalpells etwas eindrückt. Diese Deformation direkt mit den Volumendaten zu simulieren, ist sehr aufwändig, und nicht interaktiv möglich, da große Teile des eher langsamen Volumenspeichers immer wieder neu beschrieben werden müssten. Alternativ, kann ein ähnlicher Effekt erzielt werden, wenn die Sehstrahlen, wie in Abbildung 8.13 rechts, genau entgegengesetzt gekrümmt werden.

Die Krümmung kann als Versatz des Strahls von seinem ursprünglichen Verlauf angegeben werden. Die Länge und Richtung des Versatzes kann beispielsweise mit dem ChainMail-Algorithmus [83] berechnet werden, der die Kräfte innerhalb des Gewebes simulieren kann.

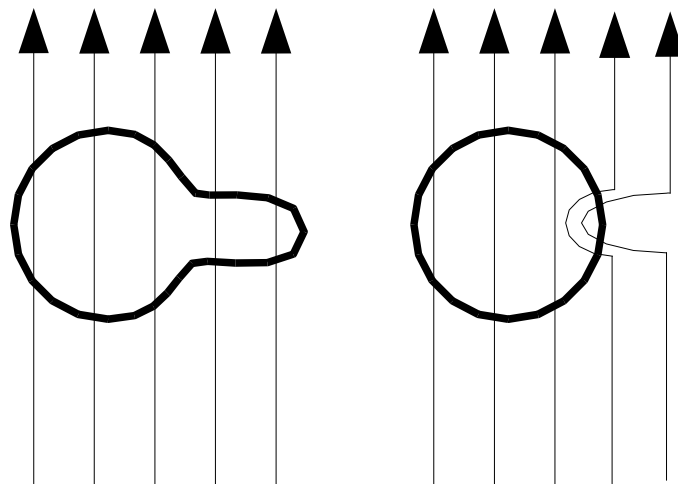


Abbildung 8.13: Deformation durch Ablenkung des Sehstrahls

Das hier vorgestellte Volume-Rendering-System kann durch die Implementierung von Space-Leaping bereits entlang des Strahls vom regulären Verlauf abweichen, in dem der Adressgenerator die Sprungdistanz berücksichtigt. Um Deformation zu unterstützen, muss der Adressgenerator noch einen Offset mit zwei Komponenten ($\Delta x, \Delta y$) senkrecht zum normalen Strahlverlauf addieren, wie in Abbildung 8.14 angedeutet.

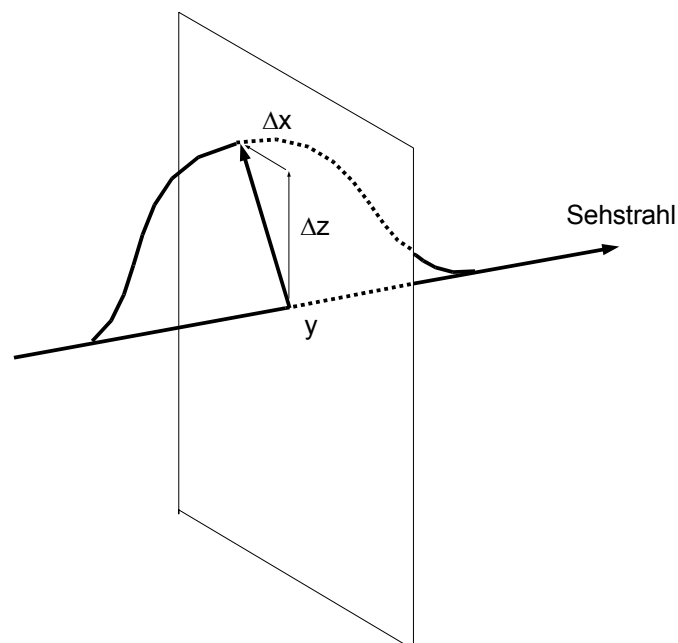


Abbildung 8.14: Offset rechtwinklig zum Strahlverlauf

Dies ist realisierbar mit einem Deformationsspeicher, der für alle 64 Strahlen, für jeden Punkt y , die

Komponenten Δx und Δy enthält. Er wird mit der Länge y und der Strahlnummer adressiert. Wird Δx und Δy mit je 2 Byte ausgelegt und die Länge des Strahls, oder deformierbaren Bereichs auf 1024 beschränkt, hat der Speicher die Größe $64 \times 1024 \times 4 \text{ Byte}$, was 256kByte entspricht. Unter den Bedingungen wie in Kapitel 7.6.3 festgelegt, mit einem 133MHz, 64 Bit PCI-Bus, ist die Übertragungszeit etwa 0.5 ms um den Speicher komplett zu füllen. Da der Speicher für jedes Strahlbündel aktualisiert werden muss, kann dies die Bildwiederholrate zwar stark reduzieren, falls große Teile des Volumens deformiert werden sollten^I. Im Normalfall wird sich die Deformation aber eher auf wenige Volumenbereiche beschränken, und nicht die vollen Länge des Strahls ausnutzen, so dass nur bei wenigen Strahlbündeln der komplette Speicher beschrieben werden muss. Bei den aktuellen Xilinx FPGAs mit mehr als 1 MByte internem Block-RAM, könnte der Deformationsspeicher intern implementiert werden.

I Bei einem 512²-Bild mit Deformation in jedem Strahlbündel summieren sich die Ladezeiten des Deformationsspeichers auf maximal 2 Sekunden.

9 Ergebnisse und Diskussion

9.1 C++-Simulation

9.1.1 Übersicht

Die Simulation der Architektur in der Programmiersprache C++ war die Grundlage dieser Arbeit. Die Architektur wie sie in Kapitel 7 beschrieben wurde und damit das Simulationsmodell, stand nicht von vornherein fest, sondern wurde mit Hilfe der C++-Simulation entwickelt. Begonnen wurde mit dem reinen Ray-Casting-Algorithmus ohne Optimierungstechniken, dem eine immer genauere Simulation realer Speicherbausteine zugefügt wurde. Hierdurch konnte ein genaues Verständnis des Zusammenspiels zwischen Algorithmus und Speicherzugriffen gewonnen werden. Nach dem Hinzufügen der Optimierungstechniken wurde deutlich, dass mit den bisherigen Verfahren keine ausreichende Leistung zu erreichen war und der Ablauf des Algorithmus sich stark an den Eigenheiten der Speicherbausteine orientieren muss.

So wurden mit der C++-Simulation neue Ideen ausgetestet und die Realisierbarkeit mit der VHDL-Simulation und Synthese nachgewiesen. Umgekehrt wurden aber auch Einschränkungen der Hardware, wie die Rechengenauigkeit und das Zeitverhalten, immer wieder in die C++Simulation eingebracht um ihren Einfluss auf die Berechnungsdauer und die Bildqualität zu überprüfen.

Folgende Funktionen wurden durch die Simulation überprüft:

- Berechnung von kompletten Ergebnisbildern zur Abschätzung der Bildqualität, wobei die Rechengenauigkeit der Hardware in die Simulation übernommen wurde.
- Überprüfung der algorithmischen Optimierungen und Berechnung der prozentualen Gewinne, bei unterschiedlichen Datensätzen.
- Einfluss unterschiedlicher Speicherorganisationen auf die Bildwiederholrate.

Dabei wurde vor allem berücksichtigt:

- Die Abarbeitungsreihenfolge der Strahlen durch Multithreading mit unterschiedlichen Sortierungskriterien
- Die Zeilenwechsel bei dynamischen Speichern.

Untersucht wurde die Abhängigkeit der Zeilenwechsel in den Speicherbänken und Bildwiederholrate von

- der Anzahl der parallel berechneten Strahlen,
- dem Rotationswinkel,

- der Klassifizierung,
- und Datensätze unterschiedlicher Auflösung.

Eine Visualisierung der Zeilenwechsel im Bild machte den Einfluss der Sortierungskriterien

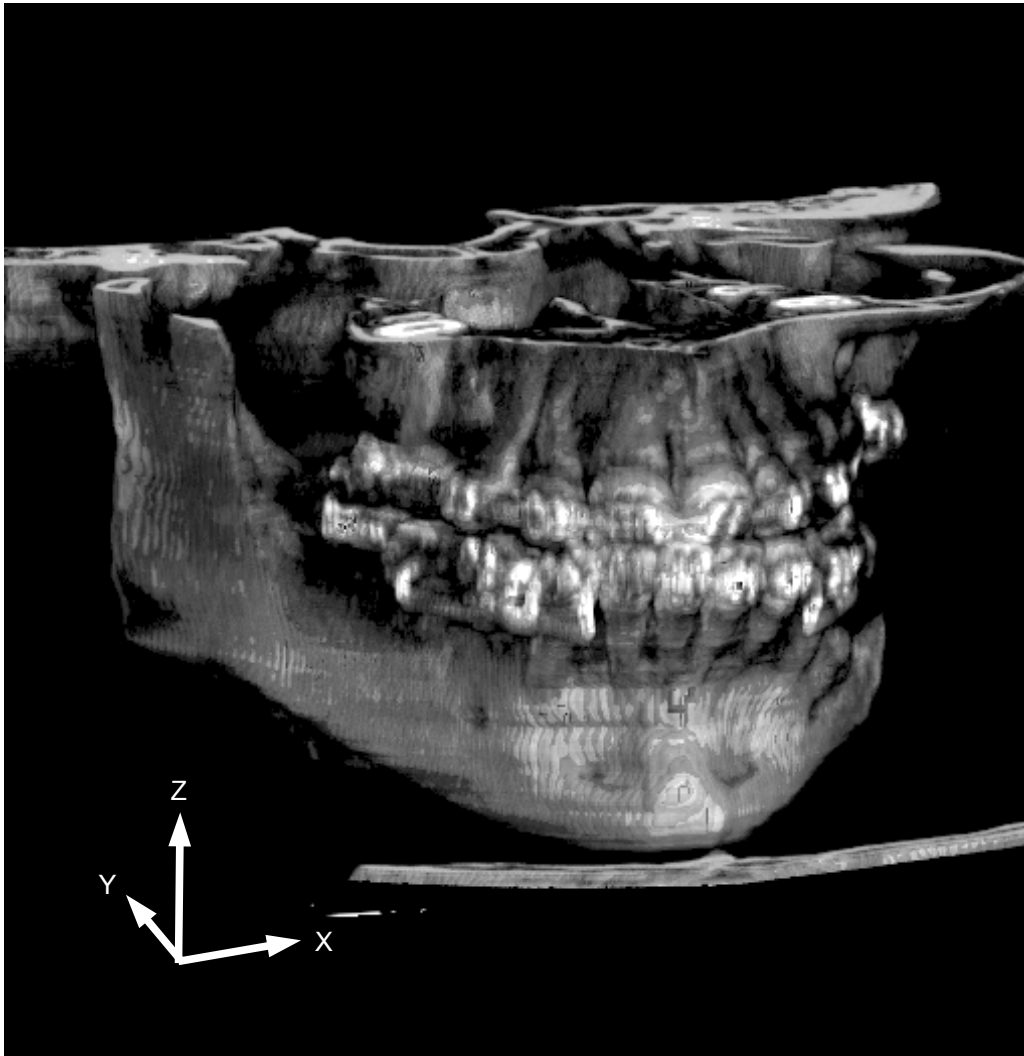


Abbildung 9.1: Menschlicher Kiefer mit Rotation von 20° um die Z-Achse und 10° um die X-Achse mit Parallelprojektion

deutlich.

Als Untersuchungsobjekt wurden Datensätze unterschiedlicher Auflösung eines menschlichen Kiefers verwendet (Abbildung 9.1). Um die Ergebnisse in den einzelnen Untersuchungen vergleichbar zu machen wurden unterschiedliche Optimierungsmöglichkeiten mit den gleichen Datensätzen, aber mit stark unterschiedlicher Klassifikation und unter Rotation des Objektes, simuliert.

9.1.2 Untersuchung der Strahlbündelgröße

Mit der in Kapitel 7.5 beschriebenen Speicherarchitektur, der eine würfelförmige Anordnung der Volumendaten im Speicherbereich zugrunde liegt, wurde der Einfluss der Strahlbündelgröße auf die Auslastung der Pipeline untersucht. Die Strahlbündelgröße entspricht der Anzahl der durch

Multithreading parallel abzuarbeitenden Strahlen. Unter Auslastung wird der prozentuale Anteil der Takte verstanden, bei denen die Pipeline einen Wert generiert. Die Takte bei denen kein Wert generiert wird, ist bei der Simulation identisch mit den Takten, die für den Zeilenwechsel der SDRAM-Speicher verloren gehen.

Bei geringer Anzahl der Strahlen müssen gleiche Speicherbereiche bei nebeneinander liegenden Strahlbündel immer wieder neu geladen werden.

Bei großen Ausdehnungen der Strahlbündel werden räumlich auseinander liegende Strahlen Daten aus verschiedenen SDRAM-Zeilen benötigen, was wiederum die Anzahl der Zeilenwechsel erhöht.

Abbildung 9.2 zeigt das Ergebnis der Simulation. Man erkennt, dass bei einer Strahlbündelgröße von 8x8, also 64 Strahlen bereits eine hohe Auslastung der Pipeline von 97 Prozent erreicht wird. Das Maximum wird allerdings bei 16x16 Strahlen mit 99 Prozent erreicht. Bei noch größeren Strahlbündel fällt die Auslastung allerdings wieder leicht ab. Dies hängt damit zusammen, dass die Speicherarchitektur 16^3 -Datenwürfel in den Zeilen der SDRAM-Bausteine, verteilt über acht Module, zwischenspeichern kann. Strahlbündel, die breiter als diese Datenwürfel sind, rufen wieder vermehrt Zeilenwechsel hervor.

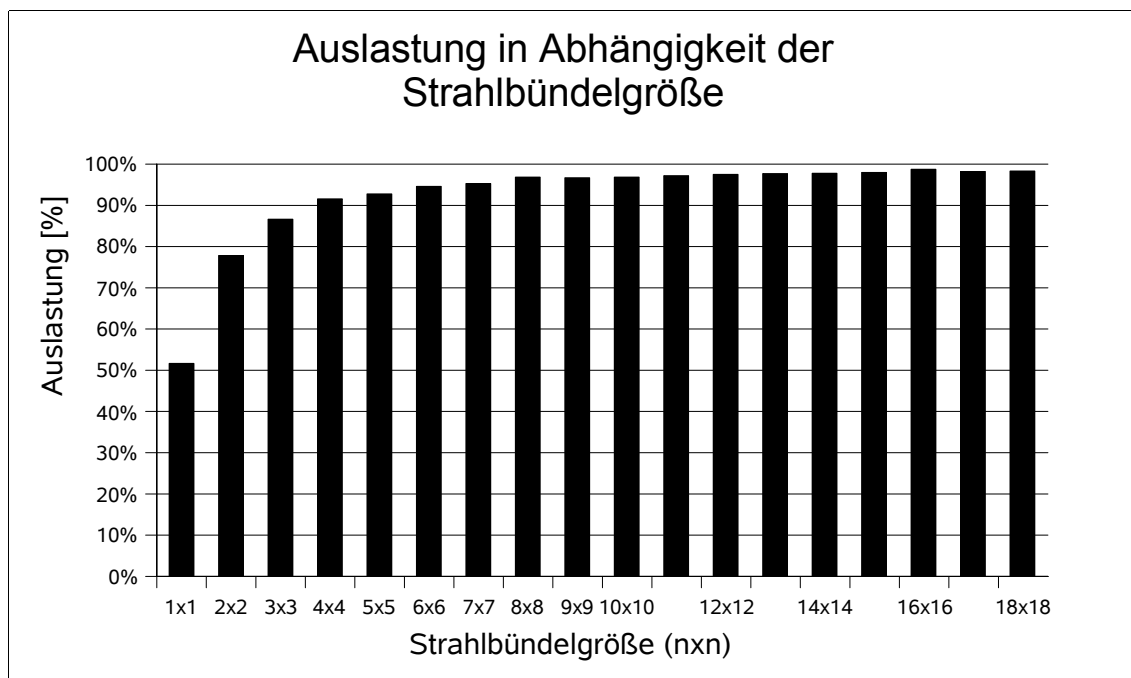


Abbildung 9.2: Pipeline-Auslastung in Abhängigkeit der Strahlbündelgröße bei einem Speichermodell mit 16^3 -Datenwürfel

Da die Verwaltung und Sortierung für 256 Strahlen bei einem 16x16 Strahlbündel sehr aufwändig wäre, wurde das Ray-Casting-System für 64 Strahlen ausgelegt und wird in dieser Form auch für die weiteren Simulationen verwendet.

9.1.3 Einfluss der Optimierungstechniken auf die Zeilenwechsel

In der folgenden Untersuchung soll der Einfluss der Optimierungstechniken Space-Leaping und Early-Ray-Termination auf die Speicherzugriffe gezeigt werden. Wie beschrieben kommt es zu Zeitverlusten, in denen die Ray-Casting-Pipeline warten muss, wenn in den SDRAM-Speichern neue Zeilen geladen werden. Das Auslesen des ersten Wertes einer neuen Zeile benötigt zehn Takte, während die folgenden in jedem Takt gelesen werden können. Bei den bisherigen Brute-Force-Verfahren ohne algorithmische Optimierungen, wie sie beispielsweise bei den Hardware-Realisierungen von Shear-Warp angewendet werden, lassen sich die Speicher regelmäßig Zeile für Zeile auslesen und so die Zeilenwechsel auf ein Minimum reduzieren. Es wird jede Zeile genau einmal geladen. Sobald aber Optimierungsverfahren verwendet werden, ist dieses regelmäßige Auslesen nicht mehr möglich und es kommt vermehrt zu Zeilenwechsel. Diese können unter Umständen so häufig vorkommen, dass ein Großteil des Gewinnes der Optimierungstechniken verloren geht.

Um dies zu visualisieren, wurden in der C++-Simulation für jeden Sehstrahl und somit für jeden Bildpunkt die Anzahl der Zeilenwechsel aufsummiert. Diese können anstatt des normalen Bildes dargestellt werden.

Zur Demonstration wird der 128x128x64 große Datensatz des menschlichen Kiefers genau von vorne, wie in Abbildung 9.3 gezeigt mit Parallelprojektion, visualisiert. Die Bildgröße ist bei allen folgenden Darstellungen 128x64 Bildpunkte, die zur Verdeutlichung vergrößert wurden. Alle Darstellungen wurden mit 8x8 Sehstrahlen visualisiert. Es wurde die in Kapitel 7.5 beschriebene Speicherarchitektur, mit acht parallelen Modulen und der Möglichkeit acht 16^3 -Datenwürfel gleichzeitig zu laden, verwendet.

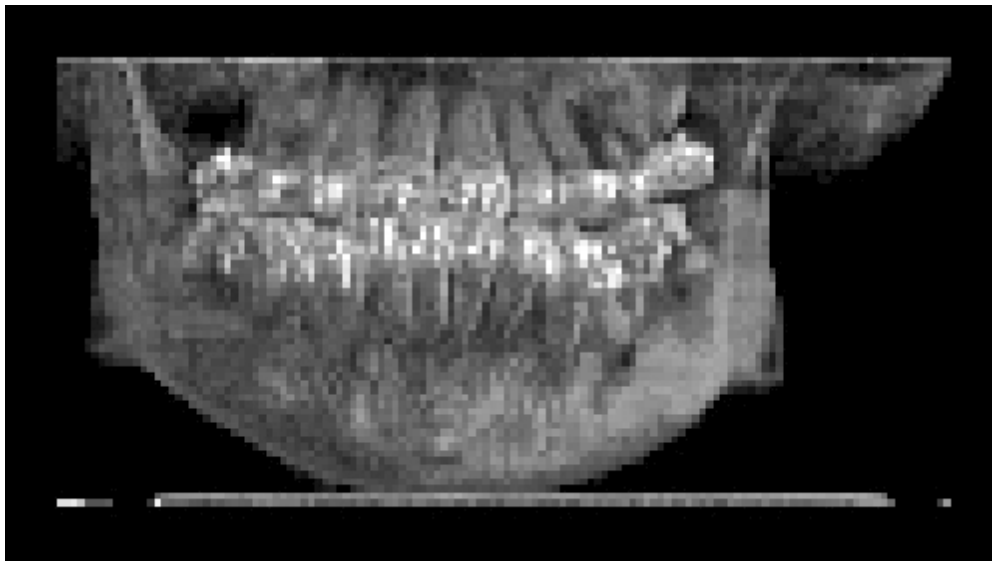


Abbildung 9.3: 128x128x64 großer Kiefer von vorne

In Abbildung 9.4 wurde der Kiefer ohne Optimierungstechniken visualisiert, weshalb ein reguläres Muster der Zeilenwechsel erzeugt wird. Für jedes 8x8-Teilbild werden nur in der rechten oberen Ecke Zeilenwechsel aufsummiert. Sie entstehen beim Übergang von einem zum nächsten 16^3 -Datenwürfel in der Tiefe des Bildes. Die Gesamtzahl der Zeilenwechsel für dieses Bild lag bei 1024. In dieser regulären Darstellung ohne Rotation des Objektes und mit parallelen Sehstrahlen, kann man dies auch analytisch nachvollziehen.

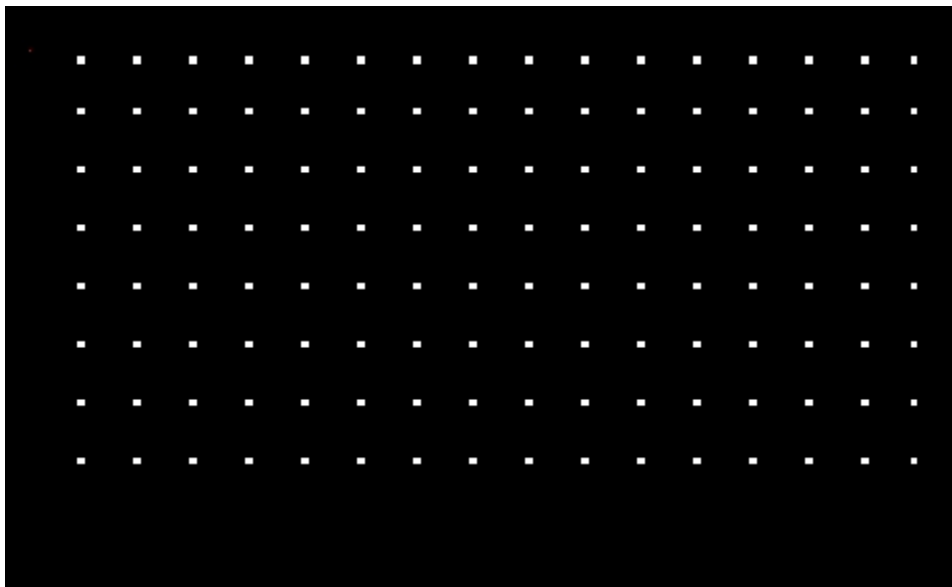


Abbildung 9.4: Darstellung der aufsummierten Zeilenwechsel pro Bildpunkt ohne Optimierungen zum Kiefer in Abbildung 9.3

Das 128x64 große Bild wird in 128 8x8-Teilbilder unterteilt. Der Datensatz ist 128

Volumenelemente tief und es kommt alle 16 Volumenelemente zu einem Zeilenwechsel. Dies entspricht 8 Zeilenwechsel pro Teilbild, multipliziert mit deren Anzahl von 128 ergibt sich die Gesamtzahl der simulierten 1024 Zeilenwechsel.

Auf die gleiche Anzahl kommt man mit der Überlegung, dass der gesamte Datensatz, von $128 \times 128 \times 64$ Größe, 256 mal 16^3 -Datenwürfel enthält. Für die 16^3 -Datenwürfel sind vier Durchgänge des 8×8 Strahlbündels erforderlich, wobei diese auch vier mal geladen werden müssen und so ebenfalls 1024 Zeilenwechsel ergeben.

Im nächsten Schritt werden bei der Visualisierung Space-Leaping und Early-Ray-Termination eingeschaltet. Die Abarbeitung der Sehstrahlen wird wieder mit 8×8 Teilbildern vorgenommen, aber immer in der gleichen Reihenfolge der Sehstrahlen innerhalb des Teilbildes von rechts oben nach links unten. Es erfolgt keine Sortierung der Sehstrahlen. In Abbildung 9.5 ist das Ergebnisbild der aufsummierten Zeilenwechsel zu sehen. Deutlich zu erkennen, ist die Häufung der Zeilenwechsel in den unteren, linken und rechten Ecken des Bildes. Ein Vergleich mit dem visualisierten Bild in Abbildung 9.3 zeigt, dass es sich um leere Volumenbereiche an den Seiten des Kiefers handelt. Hier können Teile der Sehstrahlen durch Space-Leaping schnell in das Volumen eindringen, während andere, die sich näher am Kiefer befinden, noch weiter im Vordergrund abtasten. Durch die reihum Abarbeitung der Strahlen müssen dadurch immer wieder 16^3 -Datenwürfel neu geladen werden. Die Simulation ergab 5894 Zeilenwechsel für das gesamte Bild.



Abbildung 9.5: Darstellung der aufsummierten Zeilenwechsel pro Bildpunkt mit Optimierungen und ohne Sortierung zum Kiefer in Abbildung 9.3

Der innere Teil des Bildes, in dem hauptsächlich Early-Ray-Termination zum Tragen kommt, ist keine Häufung von Zeilenwechseln zu beobachten. Hier werden die Strahlen vorzeitig beendet, so

dass weiter hinten liegende 16^3 -Datenwürfel nicht benötigt werden. Trotzdem sind noch die Eckpunkte der 8×8 -Teilbilder zu erkennen, da der Wechsel zwischen den Teilbildern immer an der gleichen Stelle über die Tiefe des Volumens erfolgt.

Für die nächste Berechnung wird eine Sortierung der Sehstrahlen nach dem Abstand zur Projektionsebene vorgenommen, so dass zuerst Strahlen bearbeitet werden, deren aktueller Abtastpunkt noch näher beim Beobachter liegt, wie in Kapitel 7.3.3 erläutert. In Abbildung 9.6 ist durch die Sortierung der Sehstrahlen keine sichtbare Abhängigkeit mehr der Zeilenwechsel vom Datensatz zu erkennen. Vielmehr sind jetzt deutlich die Ränder der 16^3 -Datenwürfel erkennbar. Die Zeilenwechsel innerhalb der 8×8 -Teilbilder sind weitgehend zufällig angeordnet, weshalb deren Struktur nicht mehr zu erkennen ist.

Der Erfolg der Sortierung lässt sich auch an der Gesamtzahl der Zeilenwechsel für diese Berechnung ablesen, die mit 1067 wieder in der Nähe des Ergebnisses ohne Optimierungen aus Abbildung 9.6 liegt. Die Sortierung verringert die Anzahl der Zeilenwechsel für diese Simulation um den Faktor 5.5.

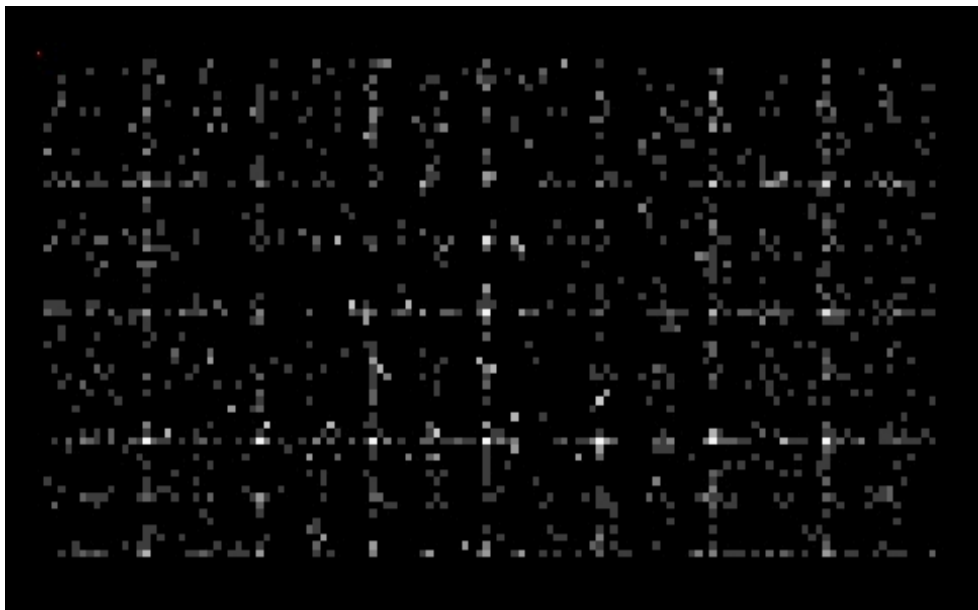


Abbildung 9.6: Darstellung der aufsummierten Zeilenwechsel pro Bildpunkt mit Optimierungen und mit Sortierung zum Kiefer in Abbildung 9.3

Bei den bisher gezeigten Berechnungen und Darstellungen verliefen die Sehstrahlen immer parallel zu den Kanten der 16^3 -Datenwürfel. Dass die Sortierung richtungsunabhängig die Anzahl der Zeilenwechsel vermindert, zeigt das Diagramm in Abbildung 9.7. Hier wurde der oben gezeigte Datensatz um die Z-Achse in 10° -Schritten gedreht und die Zeilenwechsel bezogen auf die Gesamtzahl der Speicherzugriffe jeweils mit und ohne Sortierung aufgetragen. Bei 10° und 80°

Rotation wird die Anzahl der Zeilenwechsel sogar um den Faktor 7 reduziert.

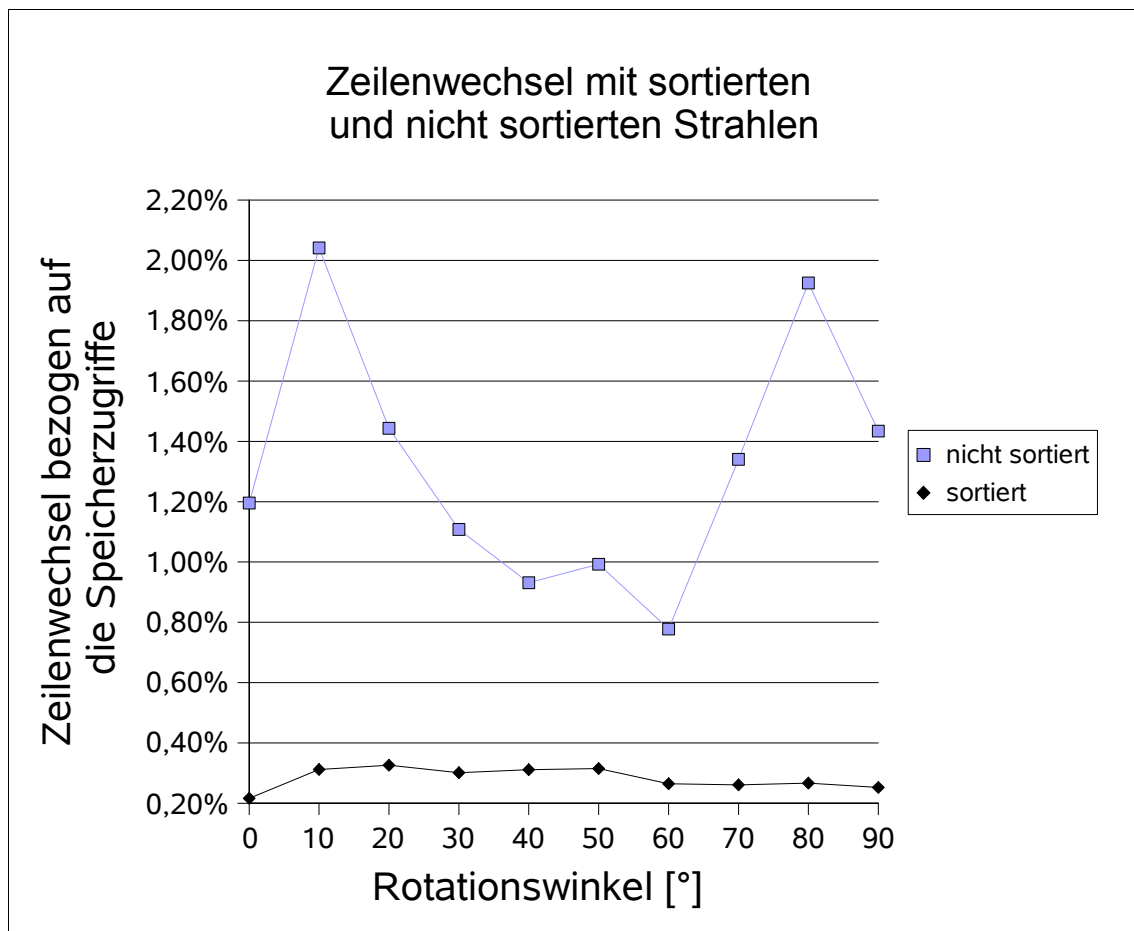


Abbildung 9.7: Prozentualer Anteil der Zeilenwechsel zu den gesamten Speicherzugriffen bei sortierten und nicht sortierten Strahlen über der Rotation um die Z-Achse

Auf die Bildwiederholrate wirkt sich die Sortierung bei diesem kleinen Datensatz um 13% aus. Wie später noch zu sehen ist, wird jede Optimierung aber um so wichtiger, je größer die Datensätze werden.

9.1.4 Unterschiedliche Sortierungskriterien

Wie in Kapitel 7.4 erläutert, gibt es mehrere Möglichkeiten, nach welchen Kriterien die Strahlen sortiert werden können. Zum einen wurde der Abstand zur Projektionsebene verwendet. Ist der Abstand gleich, handelt es sich, unabhängig von der Rotation, immer um benachbarte Abtastpunkte. Allerdings muss der Anfangswert des Abstandes für jeden Strahl vom Hostsystem berechnet werden und auf der Hardware bei jedem Abtastpunkt aktualisiert werden.

Etwas einfacher zu realisieren, ist die Sortierung nach der verbleibenden Strahllänge. Sie wird auf

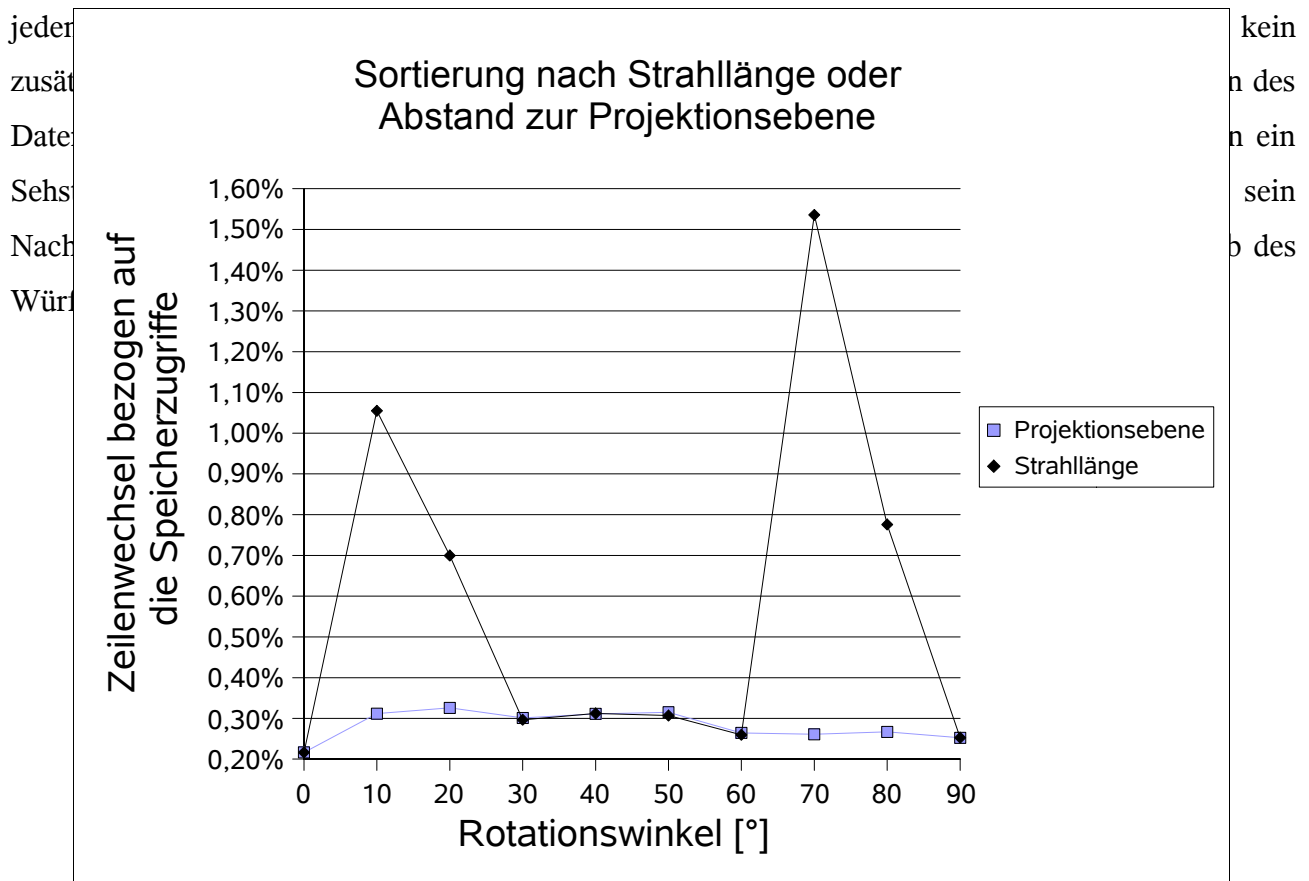


Abbildung 9.8: Zeilenwechsel bezogen auf die Gesamtzahl aller Speicherzugriffe bei unterschiedlichen Sortierungskriterien

Wie stark dieser Einfluss auf die Anzahl der Zeilenwechsel ist, wurde im Diagramm in Abbildung 9.8 untersucht. Bei den parallelen Ansichten 0° und 90° sind die Ergebnisse beider Sortierungskriterien identisch. Ebenso, wenn die Sehstrahlen mit einem steileren Winkel zwischen 30° und 60° auf den Datenwürfel treffen.

Man erkennt aber deutlich, dass bei Winkelabweichungen von 10° und 20° von den parallelen Ansichten die Zeilenwechsel stark ansteigen. Das dieser Anstieg nur am Rand des Datenwürfels erfolgt, lässt sich gut mit den visualisierten Zeilenwechsel in Abbildung 9.9 erkennen.



Abbildung 9.9: Darstellung der aufsummierten Zeilenwechsel mit Sortierung nach der Strahllänge bei 70° Z-Rotation und 0° X-Rotation

Das zugehörige Bild des Kiefers ist in Abbildung 9.10 zu sehen.

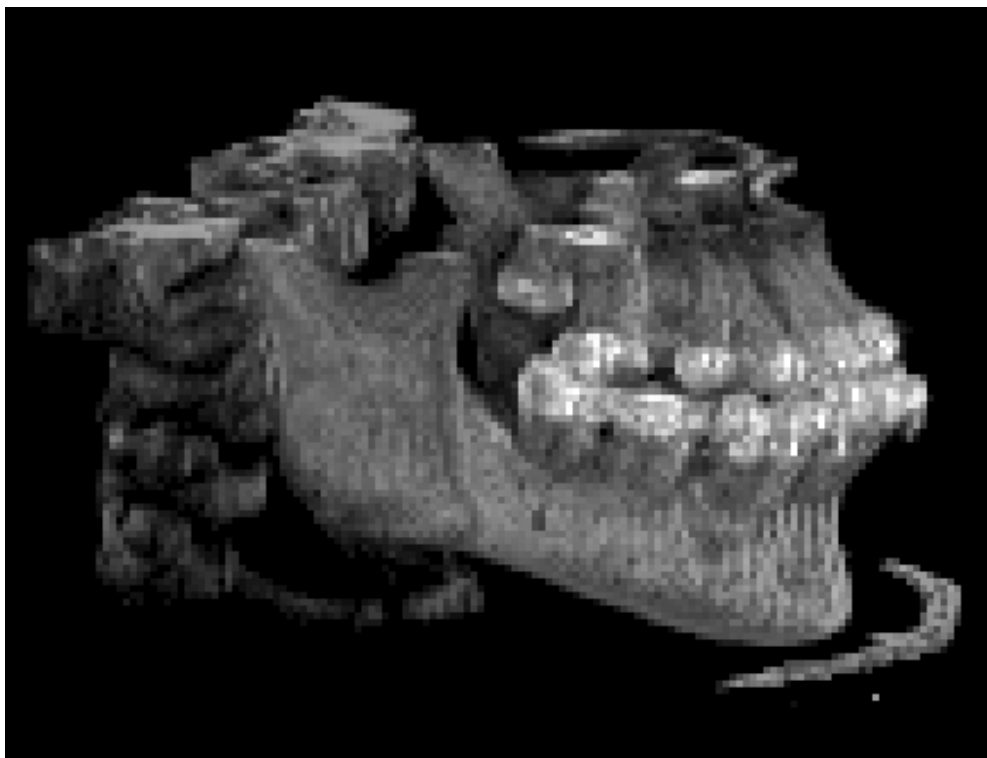


Abbildung 9.10: Kiefer unter 70° Z-Rotation und 10° X-Rotation

Der Kiefer wurde entgegen der Simulation zusätzlich noch um 10° in X-Richtung geneigt, um die Position der Steigerung der Zeilenwechsel zu verdeutlichen. Sie entstehen vermehrt an der Ecke des Datenwürfels, an der die Sehstrahlen den Würfel schneiden.

Bei der Realisierung eines Ray-Casting-Systems mit Multithreading und Sortierfunktion muss man

abwägen, ob man bei manchen Winkeln eine ein- bis zweiprozentige Erhöhung der Zeilenwechsel in Kauf nimmt oder, ob man zusätzlich die vom Host zu generierenden 64 Werte der Abstandsinformation zur Projektionsebene, bei jedem 8x8-Teilbild mit überträgt. Eine zusätzliche Rechnung auf dem Host-System ist nicht erforderlich, da der Abstand bei der Clipping-Rechnung, welche die Strahllängen bestimmt, auf jeden Fall berechnet werden muss. Auf der Ray-Casting-Hardware muss noch der Parameterspeicher um 64 10 Bit-Werte erweitert werden und ein zusätzlicher Addierer zur Aktualisierung des Abstandes mit implementiert werden, damit die Sprunginformation von Space-Leaping addiert werden kann. Auf jeden Fall bleibt die Anzahl der Zeilenwechsel weit unterhalb der Realisierung ohne Sortierung.

9.1.5 Untersuchung der Abhängigkeit verschiedener Größen bei Rotation

Bei inhomogenen Objekten werden sich die algorithmischen Optimierungen bei der Rotation unterschiedlich gut anwenden lassen, da sich die Maße der sichtbaren Oberflächen mit harten oder weicheren Kanten und Leerräume um und im Objekt ständig verändert darstellen. Hierzu wurde wieder der menschlichen Kieferdatensatz (128x128x64, Abbildung 9.1 in Parallelprojektion) verwendet. Die Rotation fand um die Senkrechte statt, so dass der Kiefer bei 90° von der Seite zu sehen ist. Hier wurde auch der Unterschied zwischen parallelen und perspektivischen Sehstrahlen untersucht. Als Perspektivenwinkel wurde 53° gewählt, wodurch sich der Abtastabstand an der Rückseite eines Würfels verdoppelt (siehe Kapitel 8.2). Daher wurde der Abtastabstand zu Beginn in der Höhe und Breite auf 0.5 gesetzt, gegenüber 1 bei der parallelen Ansicht und die Bildgröße in Breite und Höhe verdoppelt. Eine Darstellung des Kiefers in perspektivischer Projektion mit 53° Sehwinkel ist in Abbildung 9.11 zu sehen. Als Klassifikation wurden alle Werte unter 1000 ausgeblendet, wodurch nur noch die Knochen sichtbar bleiben (genauere Erläuterung im Kapitel 9.1.7 als *Linear ab 1000* beschrieben). Der Schwellwert für Early-Ray-Termination wurde auf 5% eingestellt.

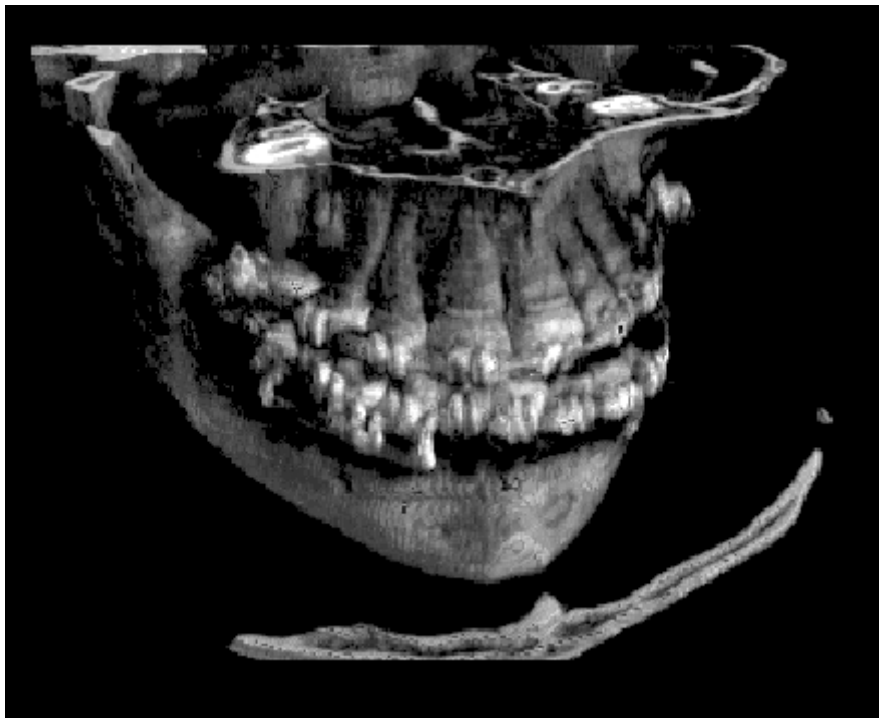


Abbildung 9.11: Menschlicher Kiefer mit Rotation von 20° um die Z-Achse und 30° um die X-Achse in perspektivischer Darstellung mit 53°

In den Abbildungen 9.12 und 9.13, sind jeweils für parallele und perspektivische Ansicht die Aufteilung des Gesamtvolumens auf den Algorithmus über dem Rotationswinkel dargestellt. Die gestreiften Anteile konnten aufgrund von Space-Leaping (17-31%) und Early-Ray Termination (34-41%) ausgelassen werden. Der graue Bereich gibt an, bei wie vielen Zugriffen nur die Distanzinformation für Space-Leaping (7-10%) gelesen wurde. Der schwarze Bereich sind die tatsächlich berechneten Abtastpunkte (21-40%), die nicht leer waren. Die ausgefüllten Teile gemeinsam enthalten somit den Gesamtanteil der Speicherzugriffe.

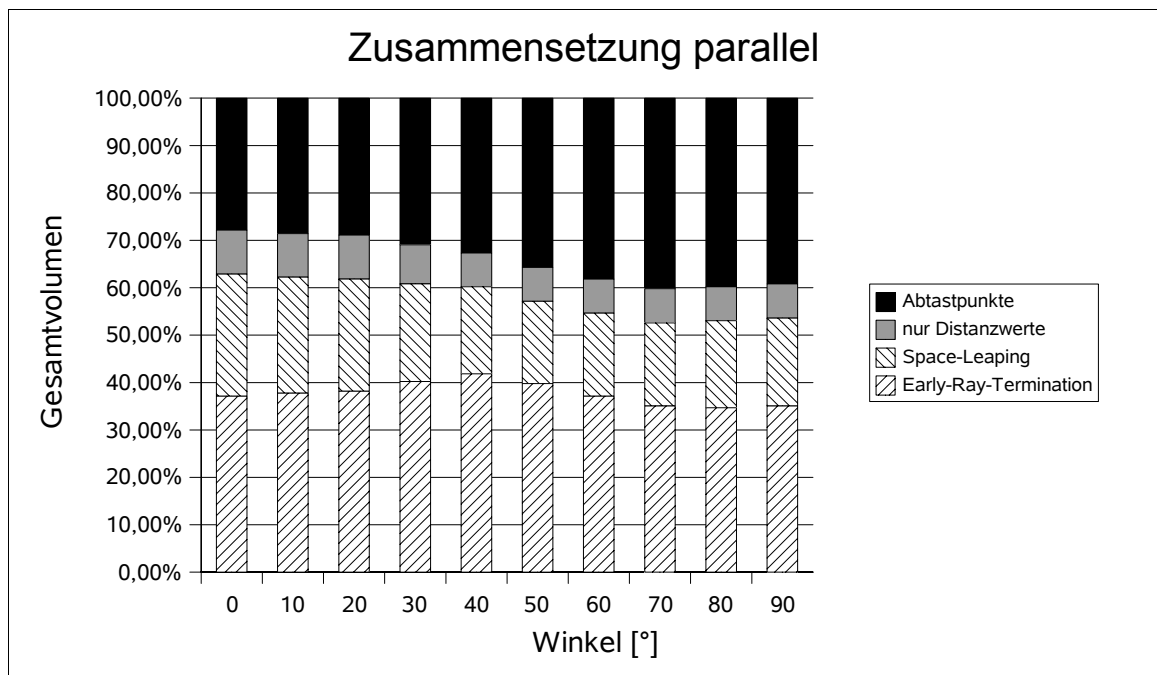


Abbildung 9.12: Prozentuale Zusammensetzung der ausgelesenen Daten und algorithmischen Optimierungen zum Gesamtvolumen mit parallelen Sehstrahlen

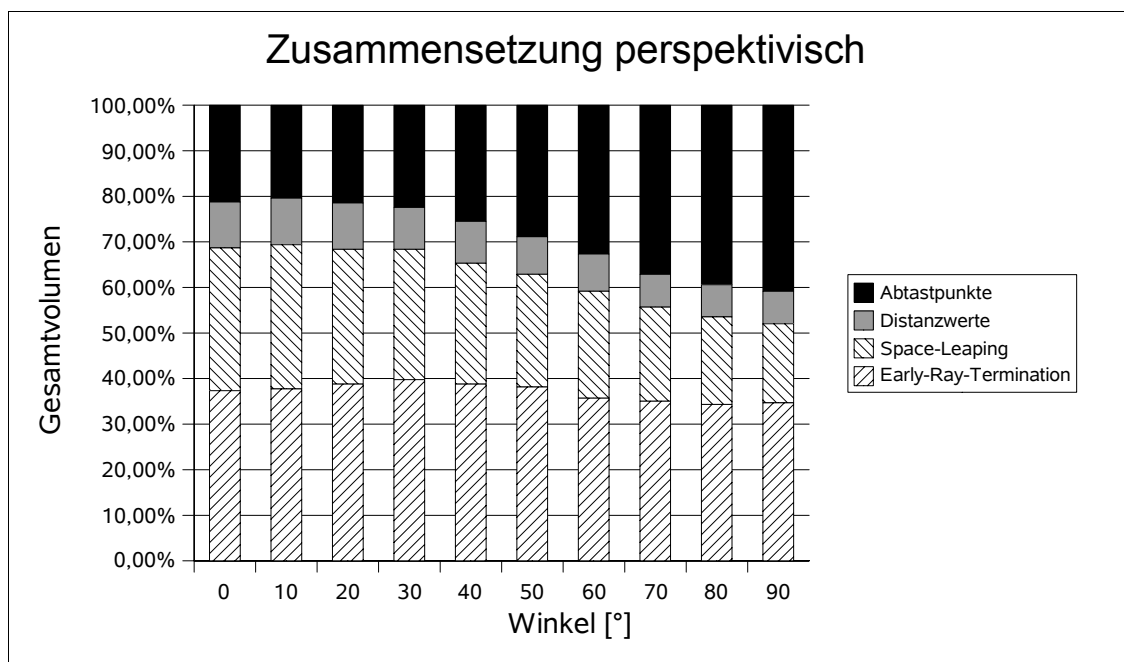


Abbildung 9.13: Prozentuale Zusammensetzung der ausgelesenen Daten und algorithmischen Optimierungen zum Gesamtvolumen mit perspektivischen Sehstrahlen

In Abbildung 9.14 ist der Anteil der Zeilenwechsel aufgetragen, der immer unter 0.37%, bezogen auf die Gesamtzahl der Speicherzugriffe, liegt. Das theoretische Minimum für SDRAM-Zeilen, die 512 Volumenelemente enthalten, liegt bei rund 0.2 Prozent ($1/512$). Dies wird erreicht, wenn jede

SDRAM-Zeile des verwendeten Datenvolumens genau einmal geladen wird. Bei paralleler Ansicht und 0° Rotationswinkel wurden diese Wert fast erreicht.

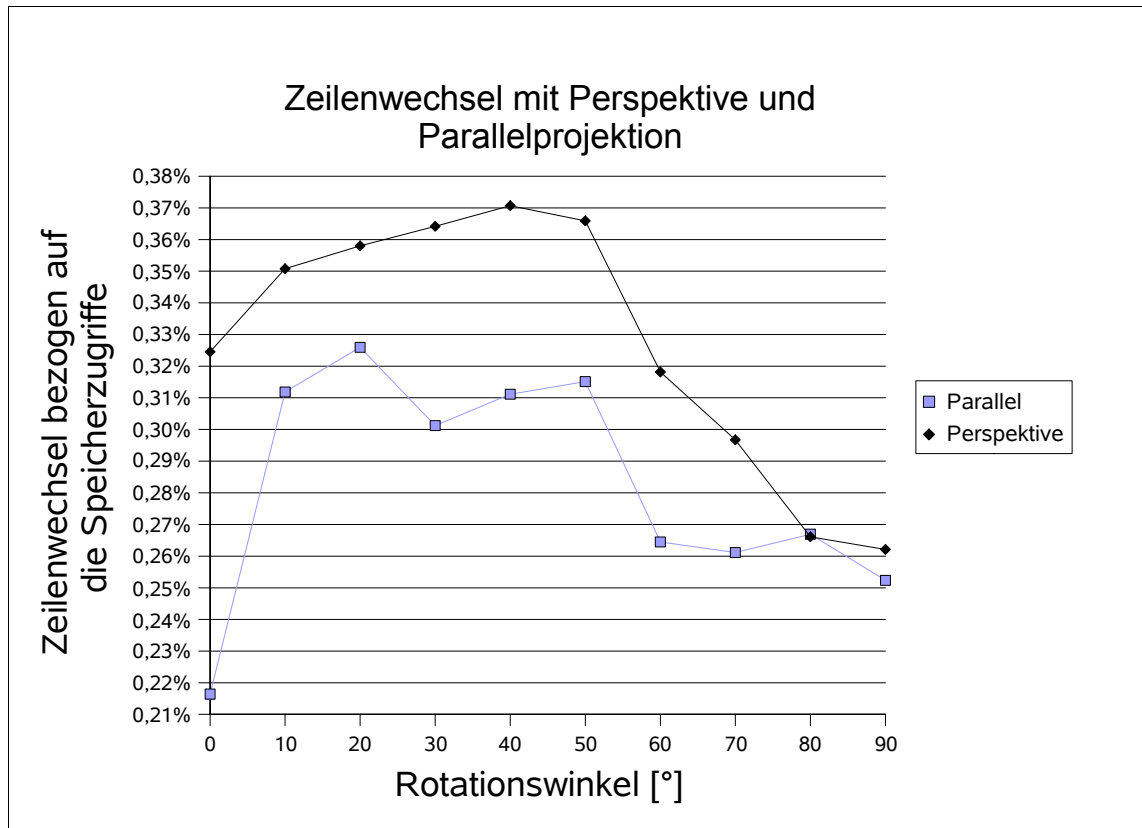


Abbildung 9.14: Zeilenwechsel mit Perspektive und Parallelprojektion in Abhängigkeit der Rotation um die Z-Achse

Der Anteil der Zeilenwechsel liegt bei Perspektive normalerweise über dem mit Parallelprojektion, da die auseinander laufenden Strahlen häufiger über die Grenzen geladener Datenwürfel gehen, als wenn sie parallel laufen.

In Abbildung 9.15 ist die tatsächliche Anzahl der Speicherzugriffe dargestellt, also die Summe aus Zugriffen auf Distanzinformation und Abtastpunkten. Bei der perspektivischen Darstellung wurde der Abtastabstand in Höhe und Breite des Bildes halbiert, und dadurch vier mal so viele Sehstrahlen durch das Volumen gelegt. Es wurde versucht den Beobachterpunkt bei Perspektive so zu legen, dass das Bild ausgefüllt ist. Trotzdem haben sich die Speicherzugriffe bei Perspektive nicht mal verdoppelt, da die Strahlen in den Randbereichen den Datenwürfel sehr bald verlassen und nicht weiter berücksichtigt werden.

Man erkennt, dass die Anzahl der notwendigen Zugriffe sich vor allem bei Perspektive während der Rotation um etwa 35% ändert, obwohl immer der ganze Kiefer sichtbar bleibt. Dies begründet sich durch den Rückgang von Early-Ray-Termination am Kiefer in der seitlichen Ansicht, da hier die

stark absorbierenden Zähne wesentlich weniger Bildfläche abdecken.

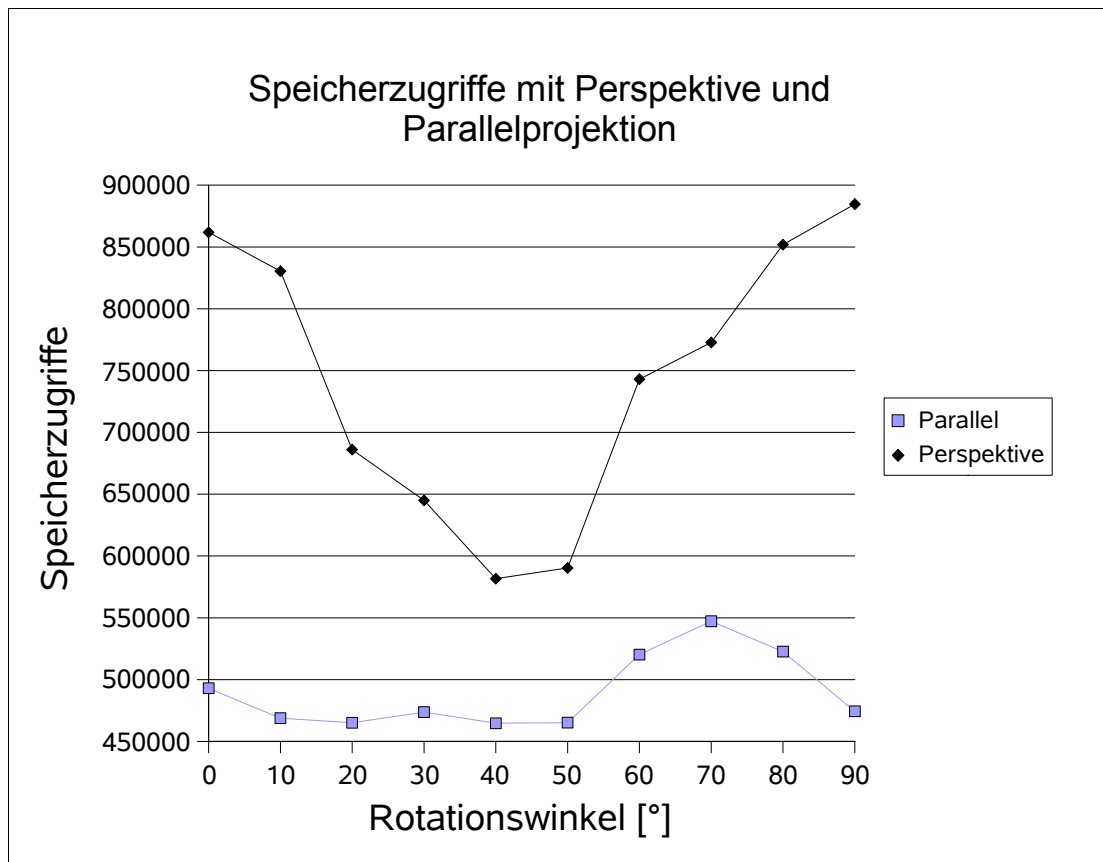


Abbildung 9.15: Abhängigkeit der Speicherzugriffe von der Rotation um die Z-Achse

Abbildung 9.16 zeigt die theoretische Bildwiederholrate eines mit 100 MHz getakteten Systems und vernachlässigt Verzögerungen, die durch Datenübertragung und Grafikaufbau entstehen könnten. Diese Verzögerungen können allerdings, durch Doppelspeicherung, Übertragung und Bildaufbau parallel zur Berechnung des nächsten Bildes, verdeckt werden. Auch hier ist zu erkennen, dass obwohl vier mal mehr Bildpunkte generiert wurden, die Bildwiederholrate bei Perspektive nur 20 bis 45 Prozent niedriger liegt.

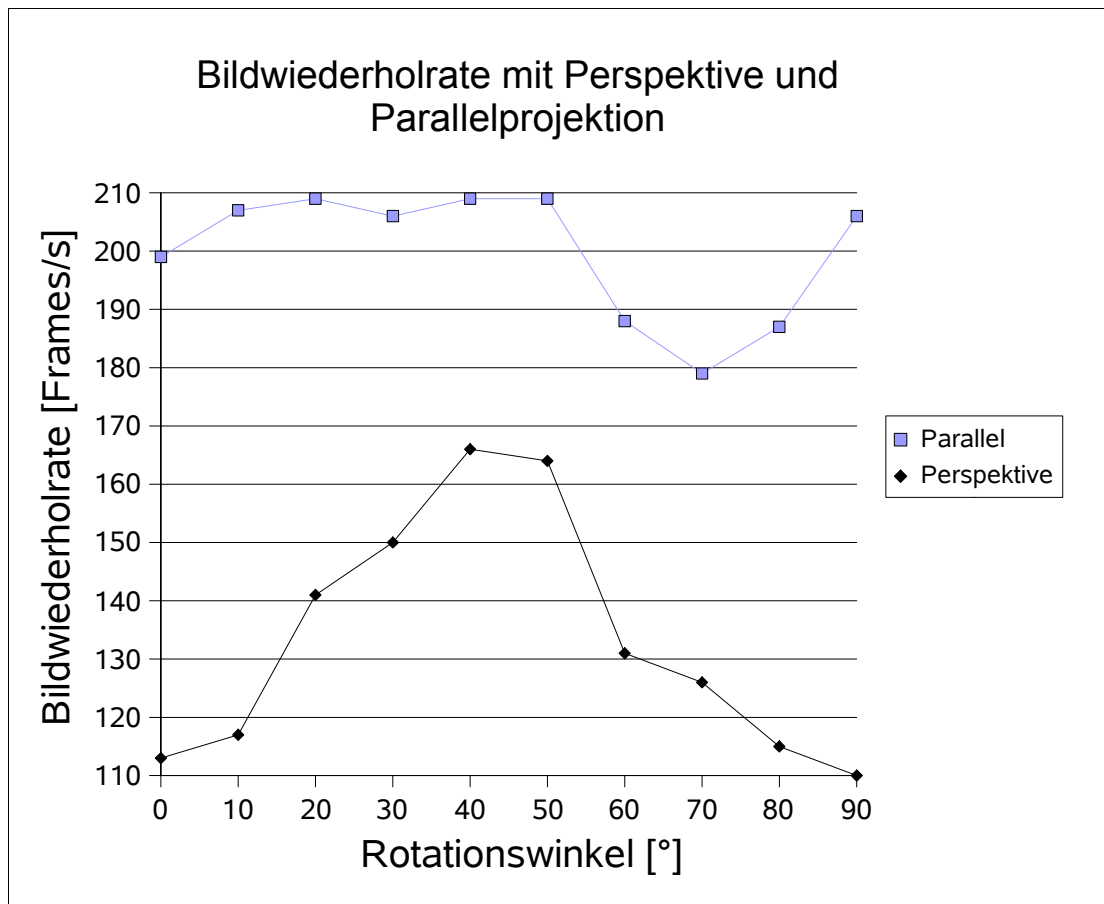


Abbildung 9.16: Bildwiederholrate in Abhängigkeit von der Rotation um die Z-Achse

9.1.6 Vergleich unterschiedlicher Speichermodelle

Hier soll gezeigt werden, wie wichtig es ist, die Volumendaten als Würfel in den SDRAM-Zeilen abzulegen, um Richtungsunabhängigkeit zu erreichen. Drei Speichermodelle werden berücksichtigt. Zum einen handelt es sich um das bereits bekannte Model mit acht Modulen und je zwei SDRAM-Bausteinen, um immer acht 16^3 -Datenwürfel laden zu können, wie es in Kapitel 7.5 beschrieben ist. Bei den beiden anderen handelt es sich zum Vergleich um Speicherarchitekturen mit nur einem Modul, bei denen für die Interpolation pro Abtastpunkt acht mal zugegriffen werden muss und somit sehr Ressourcen schonend für eine FPGA-Implementierung wäre, da nur ein SDRAM-Kontroller gebraucht wird und wesentlich weniger Anschlussleitungen notwendig sind. Der Unterschied zwischen den beiden liegt in der Zuweisung der X,Y und Z-Koordinaten zu den Adressbit. Bei einer sind die Koordinaten genauso zwischen X,Y und Z verschränkt, genau wie in Kapitel 7.5 beschrieben für die achtmodulige Speichervariante, während bei der anderen Speicherarchitektur als Negativbeispiel die Koordinaten nacheinander zugewiesen werden; die neun Bit der X-Koordinate werden den niedersten Adressbit zugewiesen, die neun Bit der Y-Koordinate den folgenden und so weiter.

Im Diagramm in Abbildung 9.17 sind die Zeilenwechsel der beiden Modelle mit verschränkten

X,Y,Z-Koordinaten über dem Rotationswinkel aufgetragen. Man erkennt, dass bezogen auf die Gesamtzahl der Speicherzugriffe, bei der Aufteilung der Daten in einem einzelnen Speichermodul

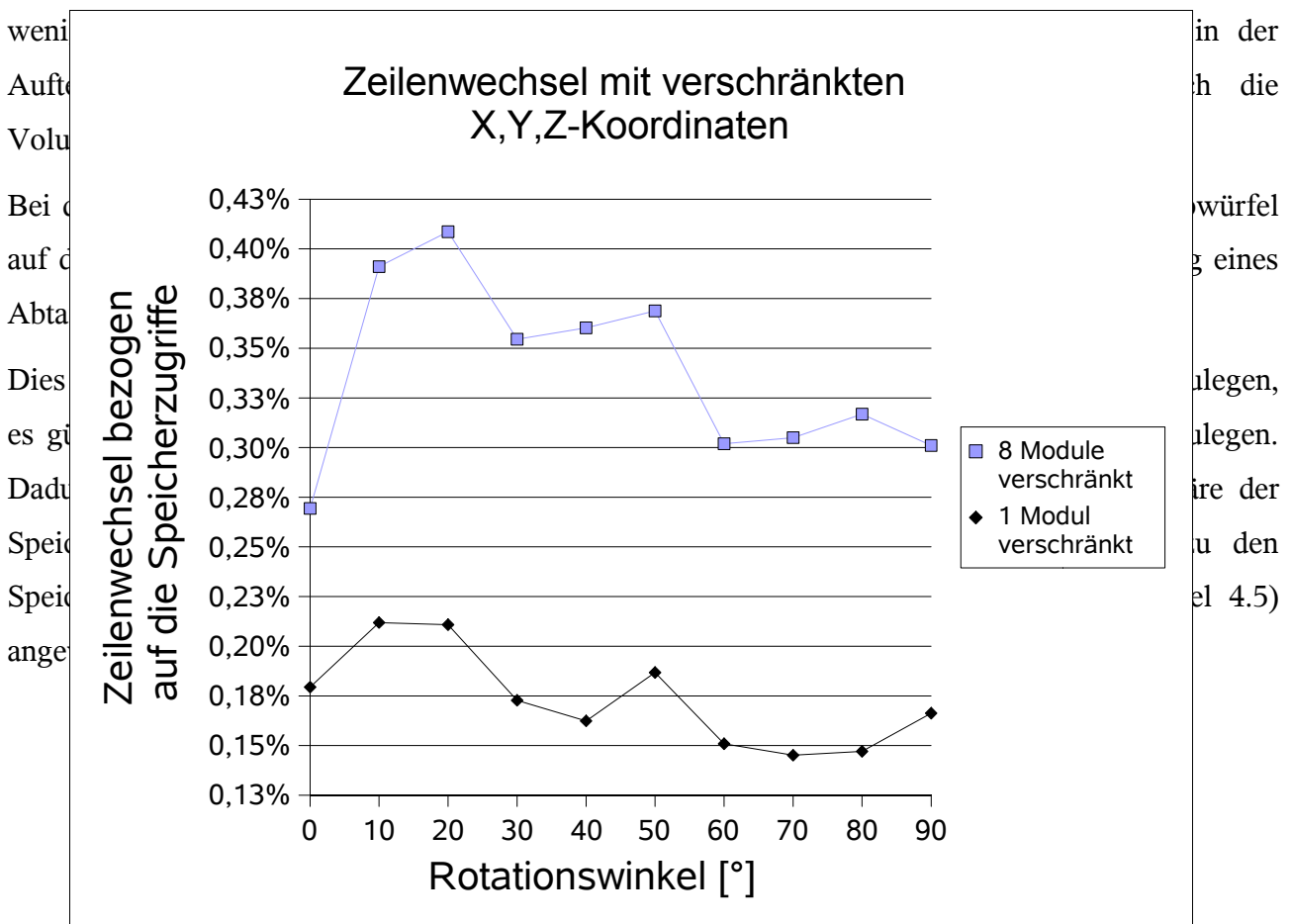


Abbildung 9.17: Zeilenwechsel mit verschränkten X,Y,Z-Adressbit bei der Rotation um die Z-Achse

Das Diagramm in Abbildung 9.18 zeigt die Zeilenwechsel der Speichervariante mit einem Modul und nicht verschränkten X,Y,Z-Koordinaten. Man erkennt, dass bei dieser Variante bei mehr als jedem zweiten Zugriff ein Speicherwechsel vorkommt. Bei der hier verwendeten Zuweisung der X,Y,Z-Koordinaten zu den Adressbit der Speicherbausteine, enthält eine SDRAM-Zeile auch eine Zeile einer Ebene des Datenvolumens. Beim Auslesen eines 2^3 -Subwürfels für einen Abtastpunkt befinden sich aber immer nur zwei Volumenelemente in einer Zeile, so dass pro Abtastpunkt mit acht Volumenelementen mindestens vier Zeilenwechsel erzeugt werden.

Genau diese Zuweisung der X,Y,Z-Koordinaten zu den Adressbit der Speicherbausteine wird erzeugt, wenn in einer höheren Programmiersprache ein dreidimensionales Feld für die Volumendaten erzeugt wird.

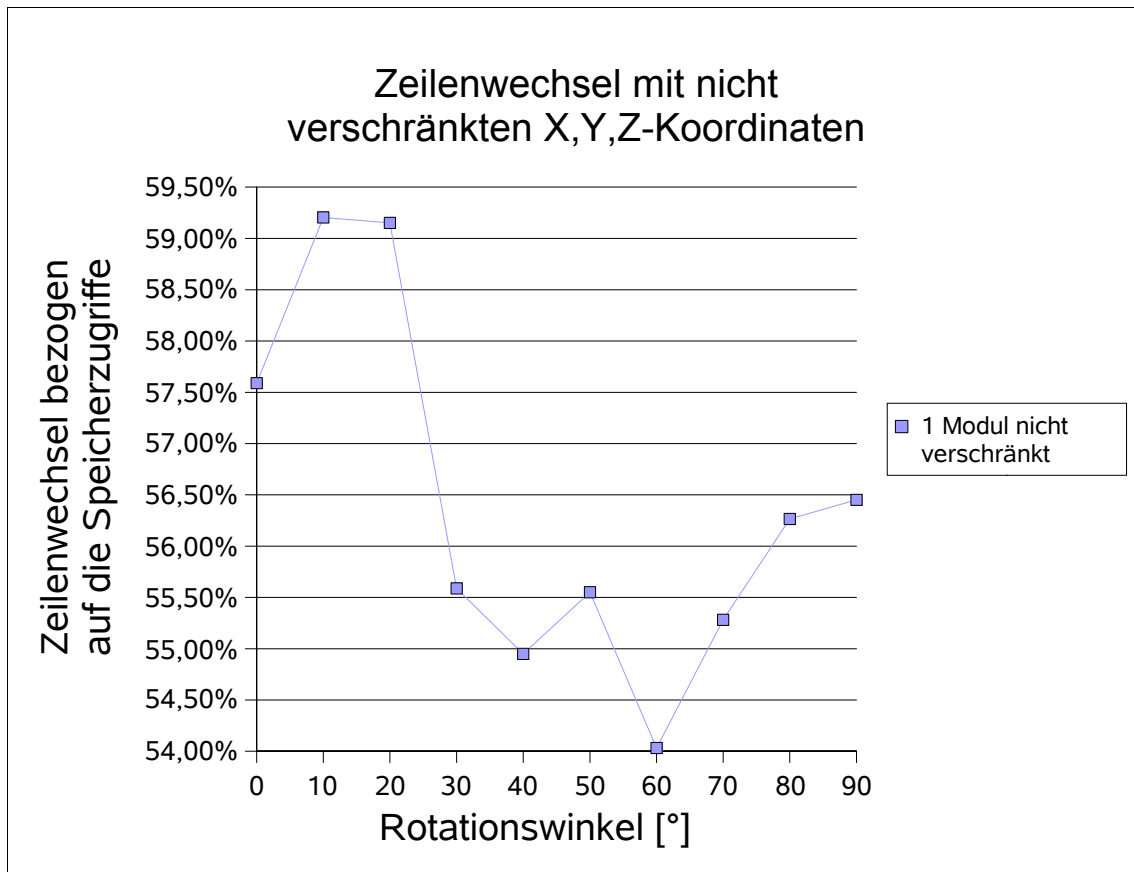


Abbildung 9.18: Zeilenwechsel bezogen auf die Gesamtzahl der Speicherzugriffe mit einem Speichermodule und nicht verschränkten X,Y,Z-Koordinaten bei der Rotation um die Z-Achse

9.1.7 Unterschiedliche Volumengrößen und Klassifizierung

Die folgende Untersuchung soll zeigen, wie sich die Bildwiederholrate bei größer werdenden Volumen und unterschiedlichen Optimierungsmöglichkeiten verhält. Hierzu wurde der Datensatz, des schon vorher verwendeten Kiefers, in einfacher (128x128x64) und doppelter Auflösung (256x256x128) verwendet und die Bildwiederholrate bei unterschiedlichen Klassifizierungen simuliert. Die Datensätze enthalten Werte von 0 bis 4095, also 12 Bit Auflösung, und müssen durch die Klassifizierung auf Werte zwischen 0 und 255 (8 Bit) umgewandelt werden, da die Ray-Casting-Pipeline auf 8 Bit ausgelegt wurde.

Erläuterung der verwendeten Klassifizierungen:

- **Schwellwert ab 1300** ist die Klassifizierung mit der größten Möglichkeit für Optimierungen. Es wurde ein Schwellwert verwendet, der allen Dichtewerten über 1300 den Wert 255 zuweist und darunter den Wert 0. Durch diese Einstellung entstehen sehr viele leere Bereiche, die durch Space-Leaping übersprungen werden können und durch die harte Schwelle werden die Strahlen mit Early-Ray-Termination schon beim ersten Kontakt mit einer Kante absorbiert und die Berechnung abgebrochen.

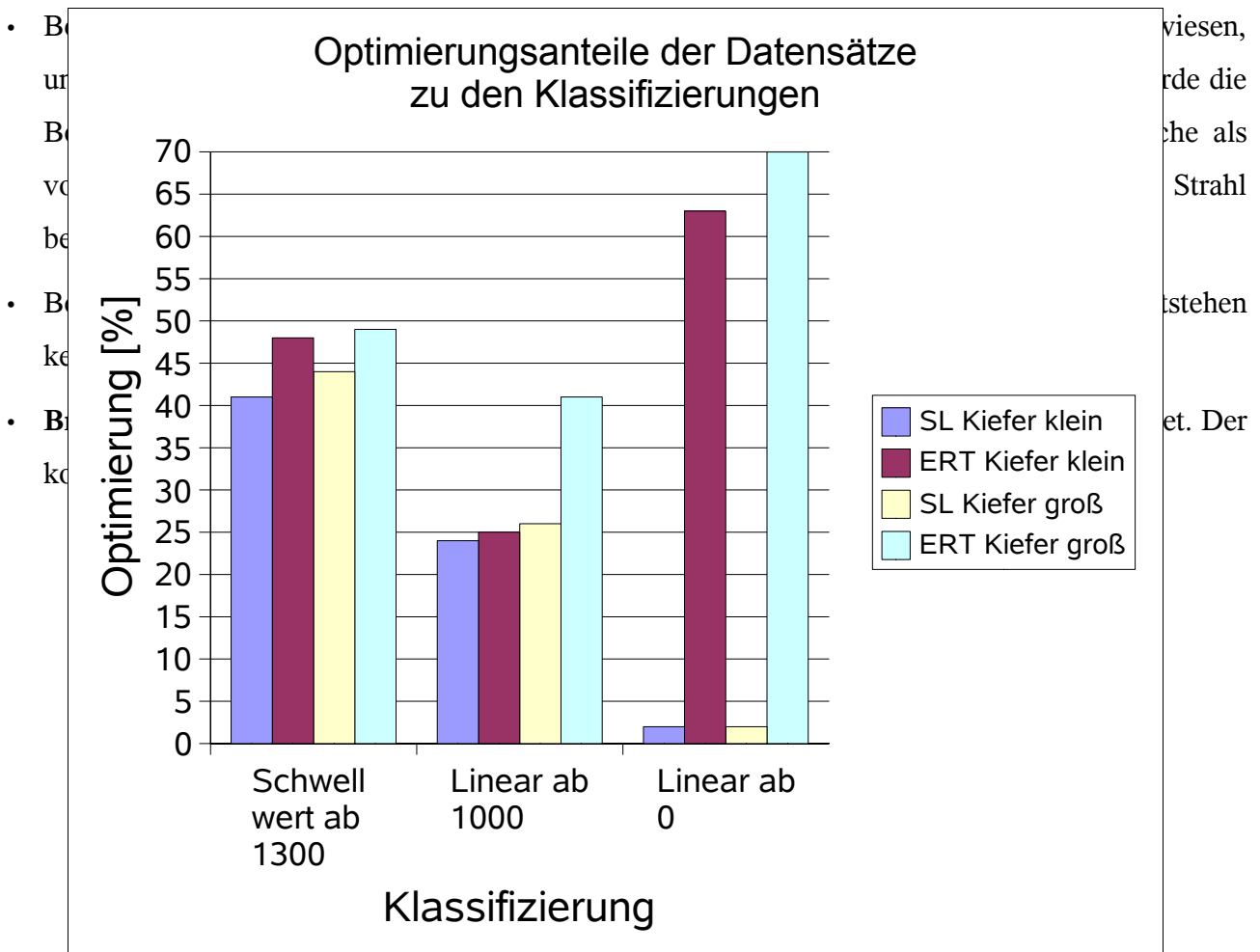


Abbildung 9.19: Optimierungsanteile der Datensätze zu den Klassifizierungen

Im Diagramm in Abbildung 9.19 sind die Optimierungsanteile der beiden Datensätze für Space-Leaping (SL) und Early-Ray-Termination (ERT) zu den verwendeten Klassifizierungen aufgetragen. Die Möglichkeiten für Space-Leaping nehmen, in der Reihe der verwendeten Klassifizierungen nach, immer weiter ab. Bei **Linear ab 0** fällt die Zunahme von Early-Ray-Termination auf. Bei dieser Klassifizierung wird die Haut des Kiefers direkt an der Grenze des Datensatzes sichtbar, so dass die Strahlen kaum eindringen können und sehr bald beendet werden. Auch bei der Bildwiederholrate macht sich der starke Einfluss des Early-Ray-Termination bei **Linear ab 0** gegenüber **Linear ab 1000** bemerkbar, obwohl man im ersten Moment erwarten würde, dass durch die zusätzlichen Leerräume bei letzterem die Bildwiederholrate höher sein müsste.

Die erreichten Bildwiederholraten sind in Tabelle 9.1 zu sehen. Bei Brute Force erkennt man, wie zu erwarten, dass bei einem achtfach größeren Datensatz, auch die Bildwiederholrate um den Faktor 8 kleiner wird. Am Faktor in der Tabelle sieht man, dass je mehr die algorithmischen Optimierungen greifen können, die Bildwiederholrate immer weniger, zum größeren Datensatz hin, abfällt und beim Schwellwert ab 1300 sogar auf 4,8 heruntergeht.

Mit Hilfe dieser Faktoren wurde die Bildwiederholrate für Datenvolumen der Größe 512^3 und 1024^3 abgeschätzt, wobei ein Datensatz mit 1024^3 Volumenelementen je nach Optimierungsmöglichkeiten mit etwa einem bis halben Bild pro Sekunde visualisiert werden kann.

Der Rechenaufwand von Volume-Rendering, der bei Vergrößerung des Datenwürfels mit der Kantenlänge kubisch ansteigt ($O(n^3)$) geht in Richtung eines quadratischen Anstiegs ($O(n^2)$). Dies verdeutlicht, wie wichtig die Implementierung algorithmischer Optimierungen wird, vor allem wenn die Datenvolumen größer werden und noch eine Visualisierung in Echtzeit möglich sein soll.

	<i>Schwellwert ab 1300</i>	<i>Linear ab 1000</i>	<i>Linear ab 0</i>	<i>Brute Force</i>
<i>128x128x64 Kiefer</i>	883 Hz	189 Hz	278 Hz	95 Hz
<i>256x256x128 Kiefer</i>	184 Hz	36 Hz	43 Hz	11,8 Hz
<i>Faktor</i>	4,8	5,25	6,47	8
<i>Abschätzung 512^3</i>	19,2 Hz	3,8 Hz	4,5 Hz	1,25 Hz
<i>Abschätzung 1024^3</i>	2,4 Hz	0,4 Hz	0,4 Hz	0,09 Hz

Tabelle 9.1: Bildwiederholrate in Abhängigkeit der Klassifizierung und Datensatzgröße

9.2 VHDL-Simulation

Im Gegensatz zur C++-Simulation, soll die Simulation in VHDL nicht die Wirksamkeit der algorithmischen Optimierungen, sondern die Realisierbarkeit in einer FPGA-Hardware oder einer integrierten Schaltung zeigen. Hierzu wurden alle Komponenten in einer synthesesfähigen VHDL-Beschreibung modelliert, die auch Register-Transfer-Level (RTL) genannt wird. Das bedeutet, dass alle getakteten Elemente der Schaltung und somit ihre Verzögerung in Taktzyklen und die kombinatorische Logik zwischen diesen Elementen beschrieben wurde.

Trotzdem sollte das Gesamtsystem möglichst nachgebildet werden, um das Zusammenspiel der einzelnen Komponenten zu überprüfen. Hierzu mussten die nicht synthesesfähigen Bauteile zumindest in ihrem Verhalten modelliert werden. Hierbei handelte es sich um die SDRAM-Bausteine des Volumen-Speichers mit den SRAM-Bausteinen für die Pre-Klassifikation. Da das System als Zusatzkarte in einem PC realisiert werden sollte, muss es durch PCI-Bus-Zugriffe gesteuert und die Ergebnisse ausgelesen werden können. Als Hardwarekomponente wurde hierzu

eine PCI-Local-Bus-Bridge modelliert. Die Simulation der PCI-BUS-Ein- und Ausgaben werden über Textdateien vorgenommen, die in Local-Bus-Zugriffe umgewandelt werden (Abbildung 9.20).

Mit dieser Simulation lassen sich zwar keine kompletten Bilder berechnen, aber die Initialisierung und Abarbeitung eines kompletten Strahlbündels zu einem 8x8 großen Teilbild überprüfen. Auch das Datenvolumen ist beschränkt, so wurde maximal eine Größe von 16^3 verwendet. Die Korrektheit der Berechnung wurde durch Vergleich der Ergebnistextdatei mit der C++-Simulation überprüft.

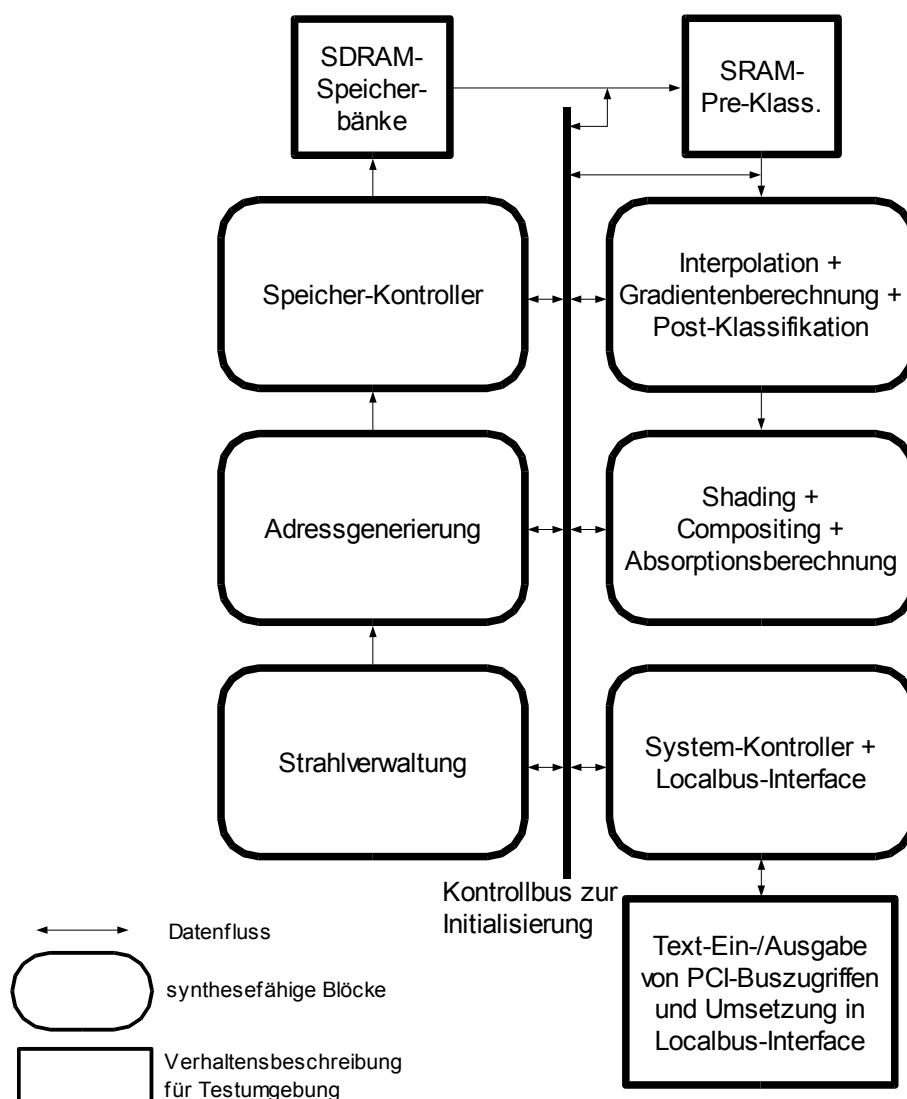


Abbildung 9.20: Aufbau der VHDL-Simulation mit Datenfluss und Unterscheidung der synthesefähigen und nicht synthesefähigen Blöcke.

Ein kompletter Simulationslauf für die Berechnung eines 8x8 Strahlbündels muss alle Initialisierungsschritte wie in Kapitel 7.6 beschrieben durchlaufen:

- Übertragen der Daten in den Volumenspeicher
- Parameter der einzelnen Strahlen in der Strahlverwaltung setzen
- Kontrollregister setzen und Rendering starten
- nach Beendigung Ergebnisbild auslesen

Abbildung 9.21 soll die Art der Ein- und Ausgabe während der VHDL-Simulation verdeutlichen. Die Eingabe wird in Form einer Textdatei PCInput.txt realisiert. Dort können Einstellungen wie Taktfrequenz und Wartezeiten zwischen den Befehlen angegeben werden (LB_CLK, WAITING) und auch die Art der Speicherzugriffe mit Adresse, Wert und einer beliebigen Identifikationsnummer.

Im Beispiel:

SWR_ACC: Single Write Access

SRD_ACC: Single Read Access

DDW_ACC: DMA on Demand Write Access

DDR_ACC: DMA on Demand Read Access

Die Ergebnisse der Lesezugriffe können in der Datei PCIoutput.txt überprüft werden. Zusätzlich kann man zur Fehlersuche jedes beliebige Signal in einer List- oder Waveform-Ausgabe im Simulator untersuchen.

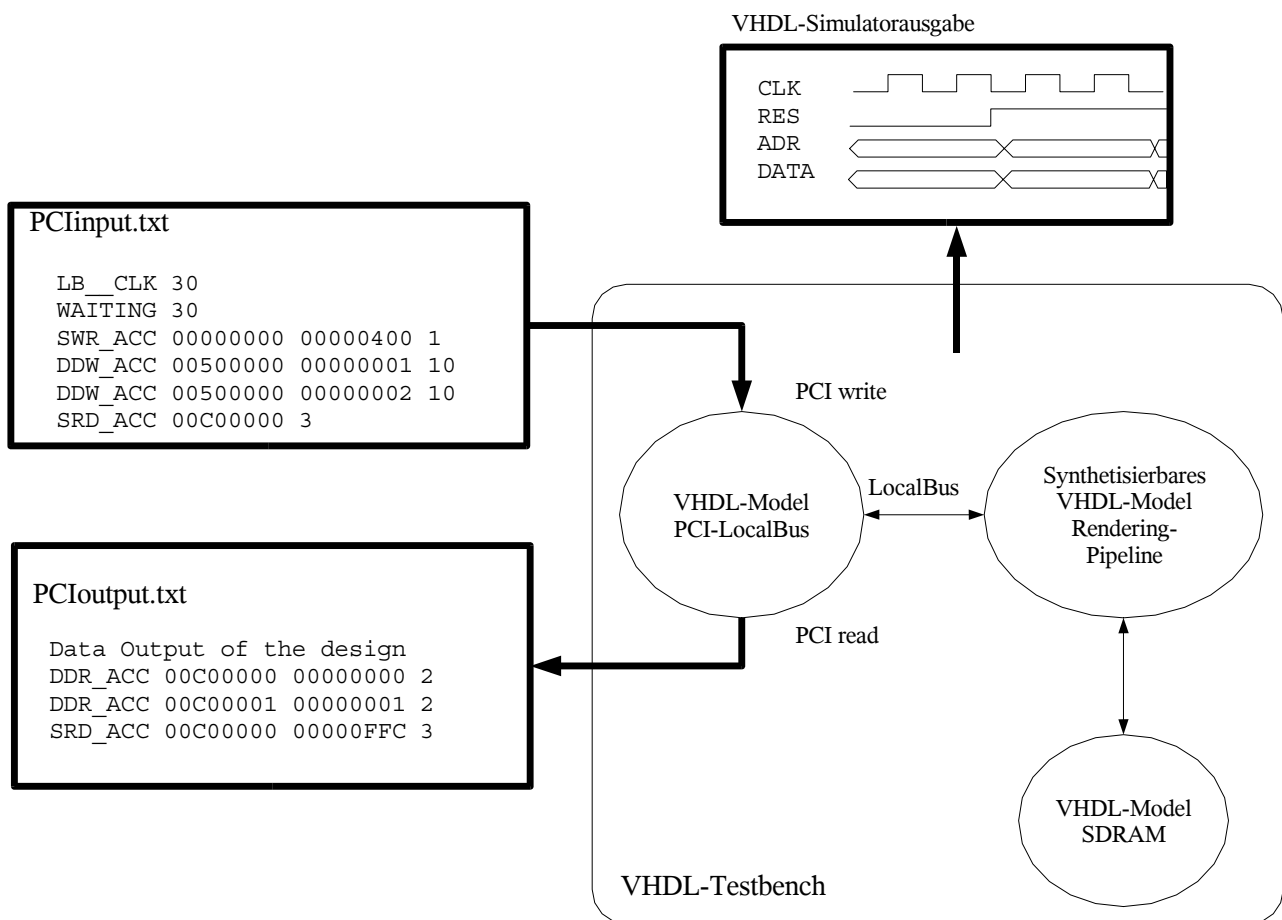


Abbildung 9.21: Ein- und Ausgabe der VHDL-Simulation

9.3 Syntheseeergebnisse

Die einzelnen Komponenten wurden zur Abschätzung der erreichbaren Taktraten und des Flächenverbrauchs für einen Xilinx-FPGA synthetisiert. Die Werte in Tabelle 9.2 wurden für einen Baustein mittlerer Größe der Virtex II-Serie ermittelt. Auf den Bausteinen stehen unterschiedliche Ressourcen zur Verfügung, die sich gegenseitig ausschließen können. So ist zum Beispiel immer ein Slice-Register mit einer 4-Input-LUT verbunden. In den 4-Input-LUTs wird die kombinatorische Logik implementiert. Wird das gebildete logische Ergebnis nicht zwischengespeichert, oder mit einem Takt synchronisiert, kann das zugehörige Register nicht mehr verwendet werden. Ebenso umgekehrt, wenn für Signalverzögerungen beispielsweise Schieberegister implementiert werden, sind die entsprechenden LUTs nicht mehr verwendbar.

Weiterhin gibt es noch spezielle 18kbit Ram-Blöcke, die allerdings über die VHDL-Beschreibung gezielt ausgewählt werden müssen. Sie können als RAM, ROM oder als FIFO konfiguriert werden.

Die 18 Bit mal 18 Bit Multiplizierblöcke (Mult-18x18s) werden vom Synthesewerkzeug anhand des Multiplikationsoperators in VHDL automatisch ausgewählt.

Der Speicher-Kontroller wurde für 8 getrennte SDRAM-Module ausgelegt, wobei die Steuerung für

einige Befehle wie Refresh und PowerDown gemeinsam verwendet werden konnte.

Der Adressgenerator und die Strahlverwaltung sind für 64 Strahlen ausgelegt. Beim Shading wurde die Post-Klassifikations-LUT und die Reflection-Map als extern angenommen. Es ist aber mittlerweile durchaus möglich, beide intern mit Hilfe der Block-RAM-Zellen zu realisieren, da neuere Bausteine über 1 MByte RAM-Blöcke zur Verfügung stellen.

Bei der Gesamtauslastung in Tabelle 9.2 erkennt man, dass der Baustein, von den Multiplizierern abgesehen, nur zu einem Fünftel ausgenutzt wird. Es ist also durchaus möglich, die Rechengenauigkeit des Systems, zum Beispiel bei der Adressgenerierung, zu erhöhen. Vor allem haben die größten Bausteine momentan eine bis zu achtfache größere Komplexität, als der hier verwendete. Die neueste Serie Virtex4 von Xilinx verfügt sogar über zusätzliche Prozessorblöcke, wodurch auch eine Berechnung der Strahlparameter auf dem FPGA möglich wäre und somit das Host-System weiter entlastet werden könnte.

Die Möglichkeiten der Logiksynthese und Werkzeuge haben sich in den letzten Jahre auch im Laufe dieses Projektes stark verbessert. So wurde beispielsweise die Adressgenerierung anfänglich auf Lucent-Orca-FPGAs implementiert. Die Place&Route-Werkzeuge liefen mehrere Stunden, um auf den Orca-Bausteinen bei 38% Ausnutzung gerade 16MHz Taktfrequenz zu erreichen. Auf aktuellen PCs benötigen die Werkzeuge für das unveränderte Design nur wenige Minuten und erreichen 100MHz bei 3% Ausnutzung. Auf einem etwa gleich großem Virtex4-Baustein (XC4V1x40) wurden sogar 130 MHz erreicht.

Trotzdem ist die Realisierung von hochgetakteten, digitalen Systemen mit FPGAs nicht trivial, wenn sehr enge Toleranzen bei der Kommunikation mit externen Bausteinen wie SDRAM oder DDR-RAM einzuhalten sind. Die Ein-, Ausgabegeschwindigkeiten und zugehörige Taktfrequenzen, sind wegen der nicht vorhersehbaren FPGA-Routing-Ressourcen häufig ein einschränkender Faktor. Bei einem mit 200MHz getaktetem DDR-RAM stehen für das Datenfenster beispielsweise nur 2,5ns zur Verfügung, wogegen die Laufzeiten von Register zu Ausgangs-Pin im FPGA leicht 10ns erreichen können [81].

	<i>Slice-Register</i>	<i>4-Input-LUTs</i>	<i>Block-RAM 18kbit</i>	<i>Mult-18x18s</i>	<i>Taktrate</i>
Speicher-Kontroller	259	432	-	-	100 MHz
Adressgenerierung	855	500	6	12	100 MHz
Strahlverwaltung	1280	2428	-	-	100 MHz
Interpolation+Gradienten berchnung	752	996	-	20	100 MHz
Shading+Compositing+ Absorption+System- Kontroller+Localbus- Interface	265	317	3	4	135 MHz
Maximum für Xilinx XC2V2000	21504	21504	56	56	
Gesamtauslastung	16%	22%	16%	64%	

Tabelle 9.2: Flächenverbrauch und Taktraten für einen Xilinx Virtex II XC2V2000-Baustein

10 Ausblick

Der Schwerpunkt dieser Arbeit lag vor allem in der Realisierung der algorithmischen Optimierungen in eine Volume-Rendering-Hardware, das in dieser Effizienz bisher von keiner Hardware umgesetzt wurde. Dies ist die Grundvoraussetzung um den wachsenden Anforderungen in Bezug auf Datenmengen und Bildwiederholraten gerecht zu werden. Für die Simulation und die ersten Versuche am FPGA, mit seinen beschränkten Ressourcen, war aber eine minimale Ray-Casting-Pipeline notwendig.

Durch die immer größer werdende Komplexität der FPGAs ist nun aber Raum neben den Optimierungstechniken auch die Bildqualität zu verbessern und Artefakte zu verringern. Die wichtigsten Punkte hierbei sind sicherlich eine verbesserte Gradientenberechnung und genauere Integrationsverfahren beim Aufsummieren der Reflexionen entlang des Strahls.

Neben der Bildqualität können auch neue Möglichkeiten mit umgesetzt werden. Einige wie Farbe, Perspektive mit größeren Winkeln und Deformation wurden in Kapitel 8 schon erläutert und können relativ leicht umgesetzt werden. Um das Ray-Tracing-Verfahren, zur Erzeugung von Schatten, umzusetzen, so dass der Datensatz nur einmal gelesen werden muss und möglichst wenig Beschränkungen bei der Position der Lichtquelle bestehen, sind sicher noch einige Ideen notwendig.

Ein großes Manko bei Space-Leaping bleibt die Vorberechnung der Distanzwerte, wodurch eine interaktive Änderung der Klassifikation beeinträchtigt wird. Hier muss nach Wegen gesucht werden, beispielsweise durch die Einführung einer hierarchischen Datenstruktur, direkt beim Rendern leere und homogene Bereiche zu erkennen. Die Multithreading-Architektur macht es auf jeden Fall möglich, dabei entstehende Berechnungszeiten zu überbrücken. Dieses System sollte in jedem Takt einen interpolierten Abtastpunkt liefern können und trotz einer komplexeren Datenstruktur die Datenkohärenz bei der Abarbeitung paralleler Strahlen berücksichtigen.

Die bisher genannten Ziele und Erweiterungsmöglichkeiten müsste man direkt angehen, sobald man ein kommerzielles System aufbauen wollte. Ein wesentlich weiteres Spektrum an Forschungsmöglichkeiten eröffnet sich, wenn man zwei Beschränkungen der hier vorgestellten Architektur fallen lässt. Es wurden nämlich nur Skalarwerte auf einem regelmäßigen Gitter zugelassen. Zum Einen gibt es viele Problemstellungen die variable Gitterabstände notwendig machen, wie sie oft in der Finite-Elemente-Rechnung üblich sind. So zum Beispiel in der Strömungstechnik, bei der an Kanten die Daten wesentlich dichter gepackt sind, als an glatten Flächen, um Verwirbelungen zu berechnen und erkennen. Alle Bereiche mit der höchsten Genauigkeit zu modellieren würde die Rechenzeiten und Datenmengen extrem wachsen lassen.

Die zweite Beschränkung bezieht sich auf die Skalarwerte. Wollte man beispielsweise einen

Tornado visualisieren, hat man in jedem Volumenelement gleich mehrere Variablen, wie Temperatur, Luftdruck, Feuchtigkeit oder Windgeschwindigkeit, die zu visualisieren sind. Das Schlagwort für Verfahren, die mehrere Variablen im Raum anzeigen, ist *Feature-Extraction*. Die einzelnen Variablen werden hier nicht nur mit unterschiedlichen Farben und Helligkeiten dargestellt, sondern es werden auch unterschiedliche Symbole, wie Pfeile für die Windrichtung, in verschiedenen Größen im Volumen mit eingeblendet. Verwirbelungen können zum Beispiel durch virtuelle Raucherzeuger, die vom Anwender platziert werden, sehr anschaulich visualisiert werden.

Auf dem Gebiet der Volumenvisualisierung wird es also noch lange nicht langweilig und neue Möglichkeiten wecken auch neue Ideen.

11 Literaturverzeichnis

- [1] Bartz D., Meißner M., "Voxel versus Polygons: A Comparative Approach for Volume Graphics", Volume Graphics Workshop 99, Swansea, März 1999 und Buchausgabe Chen, Kaufman, Yagel, Volume Graphics, Kapitel 10, Springer, 2000
- [2] Poliwoda C., "Die Entwicklung eines Parallelrechners zur Echtzeit-Volumenvisualisierung drei dimensionaler Voxelbilder", Diplomarbeit am Lehrstuhl für Informatik V, Universität Mannheim, 1994
- [3] Günther T., Poliwoda C., Reinhart C., Hesser J., Männer R., Meinzer H.-P., Baur H.-J., "VIRIM: A Massively Parallel Processor for Real-Time Volume Visualization in Medicine", W. Straßer, 9th Eurographics Workshop on Graphics Hardware, Oslo, Norway, 1994, Seite 103-108
- [4] Hesser J., "The VIRIM project : design and realization of a real-time direct volume rendering system for medical applications", Düsseldorf: VDI-Verl., 2000
- [5] Deyssenroth M., "A Shading and Compositing Processor based on Crossbar Switches and SRAM Memory", Volume Graphics Workshop 99, Swansea, März 1999 und Buchausgabe Chen, Kaufman, Yagel, Volume Graphics, Kapitel 18, Springer, 2000
- [6] Knittel G., "A PCI-based Volume Rendering Accelerator", W. Straßer, 10th Eurographics Workshop on Graphics Hardware, Maastricht, Niederlande, 1995, Seite 73-82
- [7] Knittel G., Straßer W., "A Compact Volume Rendering Accelerator", Symposium on Volume Visualization, Washington, DC, 1994, Seite 67-74
- [8] Pfister H. et al, "Cube-4 A scalable architecture for real-time volume rendering", Visualization , Volume Visualization Symposium, 1996, Seite 47-54
- [9] Osborne R., Pfister H. , Lauer H., McKenzie N. , Gibson S. , Hiatt W., Ohkami T. , "EM-Cube: An Architecture for Low-Cost Real-Time Volume Rendering", Proceedings of the SIGGRAPH / Eurographics Workshop on Graphics Hardware, p. 131-138, Los Angeles, CA, August 1997.
- [10] Pfister H., Hardenbergh J., Knittel J., Lauer H., Seiler L., "The VolumePro Real-Time Ray-Casting System", SIGGRAPH, Los Angeles, CA, 1999, Seite 251-260
- [11] TeraRecon Inc., VolumePro , <http://www.rtviz.com>, Oktober 2002
- [12] Harvey L.R., Silver D., "RACE II: A State-Of-The-Art Volume Graphics and Volume Visualization Accelerator for PC's", IEEE Visualization 1999 Late Breaking Hot Topics, San Francisco, CA, 1999, Seite 9-12
- [13] McReynolds T., Blythe D., "Programming with OpenGL: Advanced Rendering", <http://www.sgi.com/software/opengl/advanced97/notes/notes.html>, SIGGRAPH Course 1997, Oktober 2002

- [14] Barth R., "Grafikprogrammierung mit OpenGL", Addison-Wesley, 1996 (Praktische Informatik) ISBN 3-89319-975-6
- [15] Van Scheltinga J.T., Smit J., Bosma M., "Design of an On-Chip Reflectance Map", 10th EuroGraphics Workshop on Graphics Hardware, Maastricht, Niederlande, 1995, Seite 51-55
- [16] Sramek et al, "Fast Surface Rendering from Raster Data by Voxel Traversal Using Chessboard Distance", Proceedings of Visualization, Washington, DC, 1994, Seite 188-195
- [17] Lichterman, "Design of a fast voxel processor for parallel volume visualization", W. Straßer, 10th Eurographics Workshop on Graphics Hardware, Maastricht, The Netherlands, 1995, Seite 83-92
- [18] Kanus et al, "Implementations of Cube-4 on the Teramac custom computing machine", Computers & Graphics 1997, Seite 199-208.
- [19] Brady M., Jung K., Nguyen H., Nguyen T., "Two-Phase Perspective Ray Casting for Interactive Volume Navigation", Visualization, 1997, Seite 183-189.
- [20] Kreeger K., Bitter I., Dachille F., Chen B., Kaufman A., "Adaptive Perspective Ray Casting", IEEE Symposium on Volume Visualization, 1998, Seite 55-62.
- [21] Lorensen W., Cline H., "Marching cubes: A high resolution 3D surface construction algorithm", Proc. ACM SIGGRAPH Conference, 1987, 21(4):163-169
- [22] Levoy M., "Display of Surfaces from Volume Data", IEEE Computer Graphics & Appl., 8(5), 1988, 29-37
- [23] Lacroute P., Levoy M., "Fast Volume Rendering Using a Shear-Warp factorization of the Viewing Transform", Computer Graphics, Proceedings of SIGGRAPH '94, Orlando, FL, 1994, Seite 451-457
- [24] Lacroute P., "Real-Time Volume Rendering on Shared Memory Multiprocessors Using the Shear-Warp Factorization", Parallel Rendering Symposium, Atlanta, Georgia, 1995, Seite 15-22
- [25] Schlosser G., Hesser J., Zeilfelder F., Rössl C., Männer R., Nürnberger G., Seidel H.-P., "Fast Visualization by Shear-Warp on Quadratic Super-Spline Models Using Wavelet Data Decompositions", IEEE Visualization 2005, Seite 45-52
- [26] Pfister H.-P., "Towards a Scalable Architecture for Real-Time Volume Rendering", 10th Eurographics Workshop on Graphics Hardware, Maastricht, Niederlande, 1995, Seite 123-130
- [27] Dachille F., Kaufman A., "GI-cube: an architecture for volumetric global illumination and rendering", In Proceedings of the ACM SIGGRAPH/EUROGRAPHICS Workshop on Graphics Hardware (Interlaken, Switzerland, August 21 - 22, 2000), Seite 119-128

-
- [28] Frank S., Kaufman A., "Dependency Graph Scheduling in a Volumetric Ray Tracing Architecture", SIGGRAPH/Eurographics Workshop on Graphics Hardware, 2002, Seite 127-135
- [29] Chaudhry G., Li X., "A Case for the Multithread Processor Architecture", Computer Architecture News, 22(4), Sept. 1994
- [30] R. Yagel and Z. Shi, "Accelerating Volume Animation by Space-Leaping", Proc. Visualization 1993, Seite 62-84
- [31] Zuiderveld K.J., Koning A.H., Viergever M.A., "Acceleration of Ray-Casting using 3D Distance Transforms", Proceedings of Visualization in Chapel Hill, 1992, Seite 324-335
- [32] Foley, van Dam, Feiner, Hughes, "Computer Graphics: Principles and Practice", Addison Wesley, 2d. ed., 1990
- [33] Siegmann C., "Entwicklung eines Controllers zur intelligenten Verwaltung synchroner dynamischer Speicher (SDRAM Controller)", Diplomarbeit am Institut für Entwurf Integrierter Schaltkreise, Fachhochschule Mannheim, 1999
- [34] Engelmo L., "Implementierung von Algorithmen für die Volumenvisualisierung auf einem auf FPGAs basierendem System", Diplomarbeit am Institut für Entwurf Integrierter Schaltkreise, Fachhochschule Mannheim, 2000
- [35] Datenblatt SDR-SDRAM, Infineon HYB39S256400/800/160DT(L)/DC(L) 256MBit Synchronous DRAM, <http://www.infineon.com>
- [36] Geiger, Allen, Strader, "VLSI-Design Techniques for Analog and Digital Circuits", McGraw-Hill Publishing Company
- [37] Vettermann B., Hesser J., Männer R., "Solving the Hazard Problem for Algorithmically Optimized Real-Time Volume Rendering", Volume Graphics Workshop 99, Swansea, März 1999 und Buchausgabe Chen, Kaufman, Yagel, Volume Graphics, Kapitel 19, Springer 2000
- [38] Hesser J., Kugel A., Singpiel H., Vettermann B., Männer R., "Volume Rendering FPGA Hardware", XII Symposium On Integrated Circuits and Systems Design, Natal, Brasilien, 1999
- [39] Vettermann B., Hesser J., Männer R., Singpiel H., Kugel A., "Implementation of Algorithmically Optimized Volume Rendering on FPGA-Hardware", IEEE Visualization 1999 Late Breaking Hot Topics, San Francisco, CA, 1999, Seite 13-16
- [40] Hesser J., Hinkelbein C., Kornmesser K., Kuberka T., Kugel A., Männer R., Singpiel H., Vettermann B., "ATLANTIS - A Hybrid Approach Combining the Power of FPGA and RISC Processors based on Compact PCI", International Symposium on Field Programmable Gate Arrays, Monterey, CA, 1999

- [41] Chen H., Hesser J., Vettermann B., Männer R., "An Adaptive Distance-coding Based Volume Rendering Accelerator", Proceedings of the 1st International Game Technology Conference , Hongkong, 2001
- [42] Chen H., Vettermann B., Hesser J., Männer R., "Innovative computer architecture for real-time volume rendering". *Computer & Graphics*, Volume 27, Ausgabe 5, Seite 701-713, 2003.
- [43] Brosch O., Hesser J., Hinkelbein C., Kornmesser K., Kuberka T., Kugel A., Männer R., Singpiel H., Vettermann B., "ATLANTIS - A Hybrid FPGA/RISC Based Re-configurable System", Proceedings of the 7th Reconfigurable Architectures Workshop (RAW 2000), Int'l Parallel & Distrib. Proc. Symp., Cancun, Mexico ,2000, Seite 890-897
- [44] Meißner M., Kanus U., Wetekam G., Hirche J., Ehler A., Straßer W., Doggett M., Forthmann P., Proksa R., "VIZARD II: a reconfigurable interactive volume rendering system", SIGGRAPH, 2002, Seite 137-146
- [45] Doggett M., Meißner M., Kanus U., "A low-cost memory architecture for pci-based interactive ray casting", Eurographics/SIGGRAPH Workshop on Graphics Hardware, 1999, Seite 7–14
- [46] Meißner M., Kanus U., Straßer W., "VIZARD II: A PCI-Card for Real-Time Volume Rendering", Eurographics/SIGGRAPH Workshop on Graphics Hardware, 1998, Seite 61-67
- [47] Meißner M., Guthe S., Straßer W., "Interactive Lighting Models and Pre-Integration for Volume Rendering on PC Graphics Accelerators", Proc. of Graphics Interface, Calgary, Alberta, Canada, Mai 2002, Seite 209–218
- [48] Levoy M., "Efficient Ray Tracing of Volume Data", *ACM Transactions on Graphics*, Vol. 9, No. 3, Juli, 1990, Seite 245-261
- [49] Freund J., Sloan K., "Accelerated volume rendering using homogeneous region encoding", Proceedings of the conference on Visualization '97, 1997, Phoenix, Arizona, USA, 1997, Seite 191-196,
- [50] Engel K., Kraus M., Ertl T., "High-Quality Pre-Integrated Volume Rendering Using Hardware-Accelerated Pixel Shading", Proceedings Eurographics/SIGGRAPH Workshop on Graphics Hardware, 2001.
- [51] de Boer Martijn Wim, "Simulation spezieller Aspekte der nächsten Generation von Spezialrechnern für die Volumenvisualisierung", Diplomarbeit am Lehrstuhl für Informatik V, Universität Mannheim, 1996.
- [52] Bosma M., Smit J., van Scheltinga T., "Super Resolution Volume Rendering Hardware", 10th Eurographics Workshop on Graphics Hardware, Maastricht, Niederlande, 1995, Seite 117-122

- [53] Cohen-Or D., Kadosh A., Levin D., Yagel R., "Smooth Boundary Surfaces from Binary 3D Datasets", Volume Graphics Workshop 99, Swansea, März 1999 und Buchausgabe Chen, Kaufman, Yagel, Volume Graphics, Kapitel 4, Springer 2000
- [54] Kajiya J.T., von Herzen B.P., "Ray Tracing Volume Densities", Proceedings Eurographics/SIGGRAPH Workshop on Graphics Hardware, 1984, Seite 165-174
- [55] ZUIDERVELD, K. I., KONING, A. H. J., AND VIERGEVER, M.A., "Acceleration of ray-casting using 3D distance transforms", In Proceedings of SPIE, 1808, Visualization in Biomedical Computing, 1992, Seite 324- 335
- [56] D. Cohen, Z. Sheffer, "Proximity Clouds, an Acceleration Technique for 3D Grid Traversal", The Visual Computer 11:1 1994, Seite 27-38
- [57] M. Sramek and A. Kaufman. "Fast ray-tracing of rectilinear volume data using distance transforms", IEEE Transactions on Visualization and Computer Graphics, 6(3), 2000, Seite 236–25
- [58] Mueller K., Shareef N., Huang J., Crawfis R., "High-quality splatting on rectilinear grids with efficient culling of occluded voxels", IEEE Transactions on Visualization and Computer Graphics vol. 5, no. 2, 1999, Seite 116-134
- [59] Swan J.E., Mueller K., Möller T., Shareef N., Crawfis R.,x Yagel R., "An anti-aliasing technique for splatting", IEEE Visualization'97, Phoenix, 1997, Seite 197-204
- [60] Li, W., Mueller, K., and Kaufman, A. , "Empty space skipping and occlusion clipping for texture-based volume rendering", In IEEE Visualization 2003, Seiten 42–50
- [61] S. Stegmaier, M. Strengert, T. Klein, and T.Ertl, "A Simple and Flexible Volume Rendering Framework for Graphics-Hardware-based Raycasting", In Proceedings of the International Workshop on Volume Graphics 2005 , Seiten 187-195
- [62] Daqing Xue, Roger Crawfis, "Efficient Splatting Using Modern Graphics Hardware", Journal of Graphics Tools, Vol 8. No. 3, 2004, Seite 1-21
- [63] Daqing Xue, Caixia Zhang, Roger Crawfis, "iSBVR: Isosurface-aided Hardware Acceleration Techniques for Slice-Based Volume Rendering", International Workshop on Volume Graphics 2005, Seite 207-215
- [64] Klaus Mueller, Jian Huang, Naeem Shareef, and Roger Crawfis, "High-Quality Splatting on Rectilinear Grids With Efficient Culling of Occluded Voxels", IEEE Transactions on Visualization and Computer Graphics, Volume 5, Number 2, 1999, Seite 116-143
- [65] Mueller, Klaus, Torsten Moeller, Roger Crawfis, "Splatting without the Blur", Proceedings IEEE Visualization 99, Seite 363-370
- [66] Mueller, Klaus, and Roger Crawfis, "Eliminating Popping in Sheet Buffer-Based Splatting", IEEE Visualization '98, Seite 239-245

- [67] Westover L., "Footprint Evaluation for Volume Rendering", In Computer Graphics, Proceedings of SIGGRAPH 90, pages 367 376. August 1990.
- [68] Zwicker M., Pfister H., van Baar J., Gross M., "EWA Volume Splatting" , IEEE Visualization 2001
- [69] Dachille F., Kreeger K., Chen B., Bitter I., Kaufman A., "High-quality volume rendering using texture mapping hardware", Eurographics/SIGGRAPH Workshop on Graphics Hardware, 1998
- [70] Cabral B., Cam N., Foran J., "AcceleratedVolume Rendering and Tomographic Reconstruction Using Texture Mapping Hardware", In Symposium on Volume Visualization, Seite 91-98, Washington D.C., October 1994. ACM.
- [71] Cullip T. J., Neumann U., "Accelerating Volume Reconstruction with 3D Texture Hardware", Technical Report TR93- 027, University of North Carolina, Chapel Hill, 1993.
- [72] van Gelder A., Kim K., "Direct Volume Rendering with Shading via Three-Dimensional Textures", In Symposium on Volume Visualization, Seite 23-30, San Francisco, CA, Oktober 1996. ACM.
- [73] Knittel G., "TriangleCaster-extensions to 3D-texturing Units for Accelerated Volume Rendering", SIGGRAPH/Eurographics Workshop on Graphics Hardware, Los Angeles, California, Seite 25-34, August 8-9, 1999.
- [74] Wu Y., Bhatia V., Lauer H. C., Seiler L., "Shear-image order ray casting volume rendering", In Symp. on Interactive 3D Graphics, Seite 152-162, Monterey, CA, June 2003.
- [75] Knittel G., "The Ultravis System", IEEE/ACM SIGGRAPH Volume visualization and graphics symposium 2000, Oktober 2000, Seite 71- 78.
- [76] Knittel G., "Using pre-integrated transfer functions in an interactive software system for volume rendering", Eurographics 2002 Short Presentations, Seite 119-123, Eurographics Association, 2002
- [77] UltraVis Projekt, Info und Download, September 2005, <http://www.hpl.hp.com/research/ultravis/>
- [78] Roettger S., Guthe S., Weiskopf D., Ertl T., "Smart Hardware-Accelerated Volume Rendering", In Proc. Visualization Symposium '03, Seite 231-238, IEEE Computer Society Press, 2003
- [79] Krüger J., Westermann R., "Acceleration Techniques for GPU-based Volume Rendering", In Proceedings of IEEE Visualization 2003, Seite 287-292
- [80] Weiler M., Kraus M., Merz M., Ertl T.: Hardware-Based Ray Casting for Tetrahedral Meshes. In Procceedings of IEEE Visualization 2003, Seite 333-340

- [81] Vishwanathan L., Schaffer D., Tomlinson J., "Externes Gedächtnis für FPGAs", *Elektronik* 6, 2005
- [82] de Boer M., Gröpl A., Hesser J., Männer R., "Reducing Artifacts in Volume Rendering by Higher Order Integration", *Late Breaking Hot Topics Proc., IEEE Visualization* Phoenix, AZ, 1997, Seite 1-4,
- [83] Schill M.A., Gibson S.F.F., Bender H.-J., Männer R., "Biomechanical Simulation of the Vitreous Humor in the Eye Using an Enhanced ChainMail Algorithm", *Proc. MICCAI 98, Lecture Notes in Computer Science*, volume 1496, Springer, Berlin, 1998, Seite 679-687
- [84] <http://de.wikipedia.org>, 2005
- [85] Meinzer H.P., Meetz K., Scheppelmann D., Engelmann U., Baur H.J., "TheHeidelberg Raytracing Model", *IEEE - Computer Graphics and Applications*, Vol.: 11/6, 1991, Seite 34