

Universitatea POLITEHNICA București

Facultatea de Automatică și Calculatoare

**ARHITECTURI ȘI TEHNICI DE PARALELIZARE  
PENTRU SERVERE DE SPAȚII VIRTUALE 3D  
MMO**

**ARCHITECTURES AND PARALLELIZATION  
TECHNIQUES FOR 3D MMO SERVERS**

Doctorand:

Ing. Victor-Iosif Asavei

Conducător științific:

Prof.Dr.Ing. Florica Moldoveanu

București 2011



# CUPRINS

<b>LISTA FIGURILOR .....</b>	<b>7</b>
<b>LISTA TABELELOR .....</b>	<b>9</b>
<b>1 INTRODUCERE .....</b>	<b>10</b>
1.1 MOTIVAȚIE ȘI OBIECTIVE .....	11
1.2 PREZENTAREA PE SCURT A CONȚINUTULUI FIECĂRUI CAPITOL .....	12
<b>2 SPAȚII VIRTUALE 3D MMO – CONCEPTE ȘI TEHNOLOGII .....</b>	<b>15</b>
2.1 CONCEPTE GENERALE .....	15
2.2 CATEGORII DE APLICAȚII MMO .....	20
2.3 TEHNOLOGII CLIENT-SIDE .....	27
2.3.1 Banda grafică (Graphics pipeline) .....	28
2.3.1.1 Banda grafică fixă .....	28
2.3.1.2 Banda grafică programabilă .....	31
2.3.2 Organizarea și reprezentarea scenei 3D .....	35
2.3.2.1 Managementul ierarhic al scenei .....	36
2.3.2.2 Reprezentări bazate pe sortare spațială .....	36
2.3.3 Texturare .....	38
2.3.3.1 Environment mapping (maparea mediului înconjurător pe suprafețe) .....	41
2.3.3.2 Simularea detaliilor fizice ale suprafețelor .....	41
2.3.4 Sunet .....	43
2.3.4.1 Efecte audio 3D .....	44
2.3.4.2 Formate Audio .....	45
2.3.5 Animația personajelor .....	46
2.3.5.1 Modelarea personajelor .....	46
2.3.5.2 Animația personajelor .....	47
2.3.6 Sisteme de particule .....	51
2.3.7 Redări speciale .....	52
2.3.7.1 Redarea vegetației .....	52
2.3.7.2 Redarea apei .....	53
2.4 TEHNOLOGII SERVER-SIDE .....	55
2.4.1 Abordarea peer-to-peer .....	55
2.4.2 Abordarea client/server .....	55
2.4.2.1 Probleme ce trebuie tratate în realizarea proiectării sistemului client/server .....	56
2.4.2.2 Încărcarea pachetelor .....	57
2.4.2.3 Frecvența pachetelor .....	58

2.4.2.4	Ordinea pachetelor și pachetele lipsă .....	58
2.4.2.5	Clienți care încearcă să păcălească sistemul .....	59
2.5	CONCLUZII .....	60
<b>3</b>	<b>ANALIZA SOLUȚIILOR ARHITECTURALE ACTUALE PENTRU SERVERE DE SPAȚII VIRTUALE 3D MMO.....</b>	<b>61</b>
3.1	FUNCȚIONALITĂȚI DE BAZĂ .....	61
3.2	ZONARE .....	62
3.2.1	Zonare netransparentă utilizatorului .....	62
3.2.2	Zonare transparentă utilizatorului .....	63
3.3	FOLOSIREA DE INSTANȚE .....	64
3.4	ARHITECTURI TRADIȚIONALE CLIENT-SERVER .....	65
3.4.1	Descriere generală .....	65
3.4.2	interacțiunea dintre aplicațiile client și server .....	66
3.4.3	Componentele principale ale aplicației server .....	66
3.4.4	Componentele principale ale aplicației client .....	68
3.4.5	Limitări și inconveniente .....	70
3.5	ARHITECTURI P2P .....	71
3.5.1	Descriere generală .....	71
3.5.2	Limitări și inconveniente .....	73
3.6	ARHITECTURI HIBRIDE .....	74
3.7	GRID / CLOUD COMPUTING .....	75
3.8	CONCLUZII .....	76
<b>4</b>	<b>TEHNICI INOVATIVE DE PARALELISM PENTRU SPAȚII VIRTUALE 3D MMO BAZATE PE GPGPU.....</b>	<b>78</b>
4.1	DE CE GPGPU .....	78
4.1.1	Abstractizări CUDA .....	80
4.1.2	Modelul de execuție CUDA .....	87
4.1.3	Elemente de limbaj CUDA .....	90
4.1.4	Integrare CUDA cu BufferObjects OpenGL .....	92
4.1.4.1	PBO (Pixel Buffer Object) .....	92
4.1.4.2	VBO (Vertex Buffer Object) .....	94
4.1.5	Evoluție capabilități CUDA .....	94
4.2	UTILIZAREA GPGPU PENTRU REDAREA 3D CLIENT-SIDE .....	98
4.2.1	Simularea foto-realistă în timp real a scenelor 3D folosind RayTracing .....	98
4.2.1.1	Descrierea algoritmului RayTracing .....	99

4.2.1.2	<i>Calculul culorii pixelului folosind un model de iluminare</i>	102
4.2.1.3	<i>Testele de intersecție</i>	108
4.2.2	RayTracing folosind tehnici GPGPU	110
4.2.3	Adaptare algoritm RayTracing pentru GPGPU	112
4.2.4	Analiză timp de execuție	114
4.2.5	Evaluarea rezultatelor	115
4.3	UTILIZAREA GPGPU ÎN CADRUL SERVERELOR MMO	117
4.3.1	Identificare operații ce pot fi paralelizate	117
4.3.1.1	<i>Descriere generală</i>	117
4.3.1.2	<i>Logica spațiului virtual</i>	117
4.3.1.3	<i>Fizica spațiului virtual</i>	118
4.3.2	implementarea simularilor fizicii folosind GPGPU	118
4.3.2.1	<i>Volum încadratoare</i>	119
4.3.2.2	<i>Forțe implicate în simulările fizice</i>	122
4.3.2.3	<i>Algoritmi</i>	123
4.3.3	Analiză timp de execuție	126
4.4	CONCLUZII	135
<b>5</b>	<b>ARHITECTURĂ ORIGINALĂ PENTRU SERVERE 3D MMO UTILIZÂND PARALELISMUL GPGPU</b>	<b>138</b>
5.1	NECESITATEA MODIFICĂRII ARHITECTURII CLASICE	138
5.2	PREZENTAREA MODIFICĂRIILOR	140
5.3	DETALIEREA OPERAȚIILOR DE ZONARE	140
5.3.1	Optimizare ce ține cont de poziționarea spațială actuală a utilizatorilor	142
5.3.2	Optimizare ce ține cont de costul operațiilor efectuate de fiecare utilizator	143
5.3.3	Optimizare ce ține cont de tipul entităților pentru care se efectuează calcule	144
5.4	ABSTRACTIZAREA MATEMATICĂ A OPTIMIZĂRIILOR PROPUSE	144
5.5	PREZENTAREA SOLUȚIEI ARHITECTURALE COMPLETE	147
5.5.1	Modulul alocare clienți și resurse	149
5.5.2	Modulul planificare sarcini	150
5.5.2.1	<i>Creare sarcini la nivel de sistem</i>	150
5.5.2.2	<i>Planificare sarcini la nivel de GPU / CPU</i>	151
5.5.3	Modulul de procesare CUDA	152
5.6	CONCLUZII	153
<b>6</b>	<b>TESTAREA SCALABILITĂȚII SOLUȚIEI</b>	<b>155</b>
6.1	SPECIFICAȚII DE FUNCȚIONARE SERVER	155
6.2	SPECIFICAȚII DE FUNCȚIONARE CLIENT	157

6.3	COMUNICAREA ÎNTRE CLIENT ȘI SERVER.....	158
6.4	ALOCAREA RESURSELOR PENTRU CALCULUL COLIZIUNILOR .....	159
6.5	REZULTATE.....	160
6.6	CONCLUZII.....	163
<b>7</b>	<b>CONCLUZII, CONTRIBUȚII ORIGINALE ȘI EVOLUȚII VIITOARE.....</b>	<b>164</b>
7.1	CONTRIBUȚII ORIGINALE .....	165
7.2	EVOLUȚII VIITOARE.....	168
<b>8</b>	<b>LISTA LUCRĂRILOR ȘI ACTIVITĂȚILOR ȘTIINȚIFICE ALE AUTORULUI</b>	<b>170</b>
8.1	LUCRĂRI LEGATE DE TEMATICA TEZEI .....	170
8.1.1	Cărți.....	170
8.1.2	Articole.....	170
8.2	ALTE LUCRĂRI PUBLICATE ÎN PERIOADA DESFĂȘURĂRII DOCTORATULUI.....	172
8.3	ALTE ACTIVITĂȚI ȘTIINȚIFICE REALIZATE ÎN PERIOADA DESFĂȘURĂRII DOCTORATULUI .....	173
<b>9</b>	<b>BIBLIOGRAFIE.....</b>	<b>175</b>

## LISTA FIGURILOR

Figura 2-1: World of Warcraft .....	22
Figura 2-2: Lord of The Rings Online .....	24
Figura 2-3: WarHammer Online.....	26
Figura 2-4: Banda grafică fixă (clasică).....	28
Figura 2-5: Operații raster .....	31
Figura 2-6: Banda grafică programabila .....	33
Figura 3-1: Zonare fără „întreruperi” a lumii virtuale.....	63
Figura 3-2: Arhitectura clasică MMO client-server .....	65
Figura 4-1: Arhitecturi GPU și CPU moderne .....	78
Figura 4-2: Comparație între NVidia GPU și Intel CPU.....	79
Figura 4-3: Ierarhie grupuri fire de execuție .....	81
Figura 4-4: Spații de memorie.....	87
Figura 4-5: Arhitectura Fermi .....	88
Figura 4-6: Conținut SM din arhitectura Fermi .....	89
Figura 4-7: Comunicare inter-GPU înainte de CUDA 4.0.....	96
Figura 4-8: Comunicație inter-GPU Peer-To-Peer în CUDA 4.0.....	97
Figura 4-9: Spațiu de adresare unificat în CUDA 4.0.....	97
Figura 4-10: Algoritm RayTracing.....	100
Figura 4-11: Calcul componentă difuză .....	104
Figura 4-12: Calcul componentă speculară.....	105
Figura 4-13: Calcul componentă speculară folosind vectorul median.....	105
Figura 4-14: Lumină de tip spot.....	107
Figura 4-15: Scena redată folosind RayTracing.....	111
Figura 4-16: Împărțirea pe GPU / blocuri SIMD a sarcinilor .....	112
Figura 4-18: Intersecție OBB .....	121
Figura 4-19: Simularea coliziunii a 8000 de obiecte folosind 4 core-uri GPU... ..	128
Figura 4-20: Calcule de fizică pe CPU – Timpi de calcul.....	130
Figura 4-21: Calcule de fizică pe CPU – FPS.....	130
Figura 4-22: Calcule de fizică pe 1 GPU – Timpi de calcul .....	131
Figura 4-23: Calcule de fizică pe 1 GPU – FPS.....	131
Figura 4-24: Calcule de fizică pe 4 GPU – Timpi de calcul .....	132
Figura 4-25: Calcule de fizică pe 4 GPU – FPS.....	132
Figura 4-26: Comparație timpi de calcul implementării operații de fizică executate în spațiile virtuale – Partea I.....	133
Figura 4-27: Comparație timpi de calcul implementării operații de fizică executate în spațiile virtuale – Partea a II-a .....	133
Figura 5-1: Alocare statică a resurselor.....	141

Figura 5-2: Alocare dinamică a resurselor ce ține cont de poziționarea spațială a utilizatorilor.....	142
Figura 5-3: Alocare dinamică a resurselor ce ține cont de un cost asociat utilizatorilor.....	143
Figura 5-4: Alocare dinamică a resurselor ce ține cont de tipul entităților.....	144
Figura 5-5: Soluție arhitecturală cu posibilități de procesare GPGPU .....	149
Figura 6-1: Prototip spațiu virtual ce implementează arhitectura propusă .....	160
Figura 6-1: Comparație între timpii de execuție ai prototipului funcțional pentru cele trei implementari .....	162



## **LISTA TABELELOR**

Tabel 3-1: Rezultate implementare RayTracing folosind GPGPU.....	116
Tabel 3-2: Rezultate implementare calcule fizică – Partea I.....	129
Tabel 3-3: Rezultate implementare calcule fizică – Partea a II-a .....	129
Tabel 6-1: Rezultate prototip funcțional arhitectură.....	162

# 1 INTRODUCERE

Jocurile de tip Massive Multiplayer Online (MMO) precum și alte tipuri de spații virtuale cum ar fi muzee virtuale, expoziții, etc, atrag un număr din ce în ce mai mare de utilizatori. Tehnologia cât și numărul de aplicații de acest fel se află într-o continuă dezvoltare și foarte probabil spațiile virtuale 3D vor reprezenta următoarea generație de interfețe de comunicare pe Internet, înlocuind sau integrând pe cele actuale ( browsere de web, instant messengers, portaluri de socializare etc ).

Aplicațiile de tip MMO se desfășoară online având ca suport o lume virtuală persistentă la care sunt conectați și în cadrul căreia interacționează simultan sute și chiar mii de utilizatori. O lume **persistentă** este o lume virtuală online care continuă să ‘funcționeze’ chiar și atunci când utilizatorul se deconectează și nu mai participă la acțiunea din lumea virtuală. Utilizatorul se poate conecta înapoi la orice moment de timp și poate să continue de unde rămăsese în momentul deconectării.

Un utilizator al unei aplicații MMO își asumă rolul unui caracter fictiv și apoi controlează acțiunile acestuia. Caracterul fictiv reprezintă ‘avatarul’ utilizatorului. Aplicațiile de tip MMO oferă utilizatorilor posibilitatea îndeplinirii unor anumite sarcini cu obiective bine definite. Spre deosebire de aplicațiile single-user, în aplicațiile MMO utilizatorii pot îndeplini aceste sarcini fie într-un mod colaborativ împreună cu alți utilizatori sau pot alege să o facă într-un mod competitiv concurând împotriva altor jucători.

Ca exemple de tipuri de spații virtuale 3D MMO putem menționa sute de jocuri MMO, simulări sociale cum ar fi *Second Life*, aplicații de e-learning, edutainment ( educație + entertainment ), simulatoare militare pentru antrenament , etc. Spațiile virtuale au devenit foarte atractive și din punct de vedere financiar, de exemplu MMO-ul *World of Warcraft ([WOW11])* având o cifră anuală de afaceri ce depășește 1 miliard de dolari.

Din punct de vedere al proiectării, unul dintre obiectivele principale ale unei aplicații MMO 3D este acela de a crea un grad de imersiune cât mai ridicat pentru utilizatorii spațiului virtual - prin calitatea simulării (grafica 3D moderna , sunete

stereo și spațiale, simulări realiste ale fizicii, etc.), dinamica interacțiunilor (timp real) și atractivitatea scenariului.

Principalul atu al unei aplicații de tip MMO este însă dat de numărul mare de utilizatori pe care îl poate suporta simultan. Aceasta este principala caracteristică prin care se diferențiază de aplicațiile single-user / multiplayer care simulează o lume virtuală 3D.

### **1.1 MOTIVAȚIE ȘI OBIECTIVE**

Principala problema a arhitecturilor curente pentru servere de spații virtuale 3D MMO este cea a scalabilității acestora, atunci când trebuie să facă față unui număr mare de utilizatori simultan, numărul de operații ce trebuie să fie executate ajungând să fie extrem de ridicat în momentele de încărcare maximă. Aceste probleme de scalabilitate apar în mare măsură datorită necesității de a menține consistența lumii virtuale.

Marea majoritate a spațiilor virtuale 3D MMO folosesc în prezent o arhitectură de tipul Client-Server. Deși acest tip de arhitectură pune la dispoziție funcționalitățile necesare unei aplicații MMO o face cu un cost mare. Pentru a putea simula spații virtuale complexe, în cadrul arhitecturii Client-Server se folosesc un număr ridicat de echipamente pentru partea de server și în ciuda numărului lor mare, clusterelor de servere sunt limitate din punct de vedere computațional iar traficul de rețea introdus de comunicarea cu clienții nu este nici el de neglijat.

Costurile mari introduse de arhitectura Client-Server limitează practic scalabilitatea acesteia, fiind necesară împărțirea spațiului virtual în mai multe versiuni distincte ale aceleiași lumi virtuale. Deși aceasta soluție permite un anumit grad de scalabilitate al aplicațiilor MMO, ea limitează gradul de realism împiedicând interacțiunea unui număr mare de utilizatori.

De asemenea, pentru a obține un nivel de performanță în prezența unui număr mare de utilizatori, producătorii spațiilor virtuale MMO trebuie să facă față unui cost financiar semnificativ pentru a menține infrastructura sistemului.

Ținând cont de aceste limitări, obiectivele principale ale lucrării de față sunt următoarele :

- Analiza soluțiilor tehnologice / arhitecturilor actuale folosite pentru redarea spațiilor virtuale 3D atât la nivelul clientului cât și la cel al serverului
- Identificarea operațiilor ce pot fi optimizate, atât la nivelul clientului cât și la cel al serverului, și implementarea acestor optimizări folosind tehnici de paralelizare GPGPU
- Propunerea unei soluții arhitecturale originale pentru a elimina sau reduce o parte din problemele curente de scalabilitate prin :
  - o Introducerea în arhitectura serverelor a capabilităților de procesare a sarcinilor de calcul folosind tehnici GPGPU
  - o Folosirea unui mecanism de alocare dinamică a resurselor care va ține cont de o serie de factori asociați utilizatorilor spațiului virtual și operațiilor executate de aceștia
- Implementarea și testarea unui prototip funcțional al soluției arhitecturale propuse

## **1.2 PREZENTAREA PE SCURT A CONȚINUTULUI FIECĂRUI CAPITOL**

### **1. SPAȚII VIRTUALE 3D MMO – CONCEPTE ȘI TEHNOLOGII**

Acest capitol prezintă conceptele specifice spațiilor virtuale 3D MMO precum și tipurile de aplicații MMO existente în prezent pe piață. De asemenea sunt trecute în revistă pe scurt și tehnologiile folosite de servere pentru a realiza comunicarea în rețea și a asigura accesul controlat al utilizatorilor la spațiul virtual.

Un aspect foarte important atunci când vorbim de simulările 3D, este acela al creării unui grad ridicat de imersiune pentru utilizatori în spațiul virtual. Pentru a realiza acest lucru, aplicațiile MMO utilizează pentru partea de client motoare grafice care folosesc tehnologii 3D moderne. În cadrul acestui capitol sunt prezentate și o parte dintre aceste tehnologii.

## **2. ANALIZA SOLUȚIILOR ARHITECTURALE ACTUALE PENTRU SERVERE DE SPAȚII VIRTUALE 3D MMO**

Acest capitol investighează soluțiile arhitecturale utilizate în prezent în implementarea de spații virtuale 3D MMO. Este prezentată în detaliu funcționarea arhitecturii Client-Server care este cea mai folosită soluție utilizată în acest moment de aplicațiile MMO. De asemenea este trecut în revistă și modul de funcționare al arhitecturilor Peer-To-Peer și hibride.

În acest capitol sunt identificate principalele avantaje dar și lipsuri pentru fiecare din arhitecturile actuale, aceste neajunsuri apărând mai ales din cauza problemelor de scalabilitate ce sunt generate de numărul ridicat de utilizatori ce accesează în același timp spațiul virtual.

## **3. TEHNICI INOVATIVE DE PARALELISM PENTRU SPAȚII VIRTUALE 3D MMO BAZATE PE GPGPU**

Acest capitol prezintă la începutul său o descriere a tehnologiei GPGPU (General Purpose on Graphical Processing Units), de ce a fost aleasă această metodă pentru optimizarea operațiilor efectuate în spațiile virtuale 3D MMO și cum se poate implementa paralelizarea operațiilor folosind arhitectura CUDA.

În continuarea capitolului sunt identificate operațiile ce pot fi optimizate la nivelul serverului cât și al clientului, propunându-se pentru acestea soluții originale de implementare ce folosesc paralelismul atât la nivel single-GPGPU cât și la nivel multi-GPGPU.

Pentru fiecare dintre implementările realizate se prezintă rezultatele obținute cât și o analiză a acestora.

## **4. ARHITECTURĂ ORIGINALĂ PENTRU SERVERE 3D MMO UTILIZÂND PARALELISMUL GPGPU**

În acest capitol se propun soluții la nivel arhitectural pentru a reduce problemele de scalabilitate cu care se confruntă arhitecturile tradiționale de servere de spații virtuale 3D MMO. Abordările propuse încearcă obținerea unei scalabilități crescute prin introducerea de capabilități de procesare single și multi GPGPU la nivelul

serverelor de spații virtuale 3D MMO cât și a unui mecanism de alocare dinamică a resurselor necesare strâns legat de procesarea GPGPU.

În continuarea capitolului se prezintă o soluție arhitecturală completă pentru a integra ideile originale propuse anterior în cadrul arhitecturilor serverelor de spații virtuale 3D MMO. Se descriu modulele prezente la nivelul arhitecturii precum și modul lor de funcționare pentru implementarea soluțiilor propuse.

### **5. TESTAREA SCALABILITĂȚII SOLUȚIEI**

În acest capitol se prezintă rezultatele unor teste ce vizează scalabilitatea soluției arhitecturale prezentate în capitolul anterior. În acest scop a fost realizat un prototip funcțional al arhitecturii și au fost implementate atât aplicația server cât și aplicația client necesară testării.

Sunt prezentate detaliat specificațiile funcționale ale aplicațiilor server și client precum și modul în care s-a utilizat aplicația client pentru a simula un număr mare de clienți ce accesează spațiul virtual.

Capitolul se încheie prin prezentarea și analiza rezultatelor obținute, care sunt foarte încurajatoare, reprezentând astfel o primă evaluare pozitivă a soluțiilor propuse.

### **6. CONCLUZII, CONTRIBUȚII ORIGINALE ȘI EVOLUȚII VIITOARE**

În partea finală a lucrării se realizează o sinteză a activităților de cercetare, a rezultatelor obținute și se precizează direcțiile viitoare de aprofundare și dezvoltare.

De asemenea, se evidențiază contribuțiile originale ce au fost aduse de către autorul tezei.

## **2 SPAȚII VIRTUALE 3D MMO – CONCEPTE ȘI TEHNOLOGII**

Majoritatea spațiilor virtuale de tip MMO care există în prezent, sunt jocuri virtuale comerciale de tipul MMOG (massively multiplayer online game).

Așa cum spune și numele, aplicațiile MMOG permit accesul simultan al unui număr mare de jucători, în prezent acest număr variind de la sute până la zeci de mii în cadrul aceleiași lumi virtuale. Conexiunea acestora se realizează prin intermediul Internetului, aceste spații virtuale fiind practic accesibile oricui în mod gratuit sau cu costuri foarte reduse.

Aceste lumi virtuale prezintă o temă ficțională, puternic influențată de fantezie, science-fiction, legende, istorie sau uneori cinematografie.

Felul în care aceasta temă este implementată în lumea virtuală variază de la simpla sursă de inspirație pentru denumiri sau aspect al elementelor virtuale, servind doar ca tematică, până la implementarea unui scenariu narativ, pe care utilizatorul îl urmează în timp, urmărind îndeplinirea unor obiective care îi permit avansarea în cadrul jocului.

### **2.1 CONCEPTE GENERALE**

Următoarele concepte generale sunt cele mai des utilizate în cadrul jocurilor MMOG :

#### **1) Avatar**

Avatarul este reprezentarea utilizatorului în lumea virtuală, cel mai adesea sub forma unui personaj 3D sau 2D animat.

Avatarul este vizibil atât utilizatorului caruia îi aparține cât și celorlalți utilizatori.

Utilizatorul are posibilitatea de a-și controla propriul avatar într-un mod mai mult sau mai puțin avansat, elementele de baza prezente în majoritatea jocurilor fiind:

- mișcarea avatarului
- efectuarea unor acțiuni specifice, spre exemplu: saluturi, dans, etc.

- setarea stării emoționale a avatarului: vesel, trist, atent, etc.

De asemenea, există posibilitatea de a particulariza reprezentarea avatarului, prin:

- setarea unor trăsături fizice: înălțime, forma corporală, facială, culoare păr, ochi etc.
- modificarea hainelor avatarului, alegerea vocii, etc.

## **2) Personaj**

Foarte apropiat de conceptul de avatar este cel de *personaj* sau *caracter*. Acesta reprezintă utilizatorul din punct de vedere al logicii lumii virtuale respective.

Un caracter este reprezentat prin avatar și este caracterizat în general prin aspecte cum ar fi:

- clasa
- nivel: în general o valoare numerică indicând avansul caracterului pe o scara valorică
- abilități
- posesiuni : echipamente, bani (virtuali), etc
- apartenența la anumite grupuri, etc

În general, mare parte din timpul de joc este dedicată îmbunătățirii nivelului, abilităților și posesiunilor personajului.

## **3) Monștrii (Mobs)**

Aceste lumi virtuale nu sunt populate doar de avatarele jucătorilor ci și de alte categorii de entități, cum ar fi monștrii sau NPC (non player characters).

Termenul consacrat „mobs” se referă la entități cu care utilizatorul este oarecum obligat de logica jocului să le confrunte pentru a îndeplini anumite obiective sau obține anumite resurse. Spre exemplu, în momentul distrugerii unui monstru, un utilizator poate obține o anumită cantitate de aur sau obiecte virtuale pe care monstrul le „conține”. Distrugerea unor monștrii poate fi condiție de avansare în cadrul unor quest-uri, etc.

Monștrii au reprezentări dintre cele mai variate: animale, extratereștrii, personaje umanoide, etc.



De asemenea, monștrii pot avea niveluri, forță și abilități diferite, fiind mai greu sau mai ușor de distrus.

Acțiunile acestor entități sunt coordonate de către sistem, în general cu ajutorul unor elemente decizionale foarte simple.

#### **4) NPC ( Non-player characters )**

NPC sunt entități ce populează lumea virtuală și nu se află sub controlul utilizatorului.

NPC nu sunt în general inamici ce trebuie distruși ci au un rol mai variat, spre exemplu ei putând fi folosiți pentru:

implementarea anumitor dialoguri cu variante predefinite în cadrul quest-urilor sau pentru a asista utilizatorul  
a oferi anumite produse sau servicii în cadrul jocului  
însoțitori ce îl ajută pe utilizator

#### **5) Posesiuni virtuale**

Fiecare avatar poate avea anumite posesiuni virtuale, asupra cărora în general are drepturi absolute sau limitate în cadrul jocului respectiv. Spre exemplu:

- arme
- echipamente speciale
- bani
- diverse resurse
- terenuri
- etc.

El poate utiliza, modifica, vinde sau distruge posesiunile sale.

Conceptul de posesiune virtuală a dus și la crearea unor sisteme de schimb/tranzacționare a acestora, majoritatea jocurilor MMOG incluzând piețe virtuale sau sisteme de licitații („auction-house”), etc.

#### **6) Quest-uri**

Un *quest* este o suită de acțiuni pe care utilizatorul trebuie să le îndeplinească pentru a obține o recompensă în cadrul lumii virtuale.

Acțiunile pot fi eliminarea unui număr de monștri de un anumit tip, eliminarea unui monstru foarte puternic, ieșirea dintr-un labirint, cautarea și găsirea unui obiect special, asamblarea unui obiect din componente, etc.

De regulă, pe parcursul quest-urilor, utilizatorul este ghidat de către un NPC, care îi explică ce are de făcut iar uneori îi acordă ajutor.

Quest-urile pot fi individuale, la nivel de echipă sau ghildă.

## **7) Grupuri de utilizatori**

Utilizatorii se pot organiza în diverse tipuri de grupuri, pentru a colabora în cadrul jocului, spre exemplu în eliminarea monștrilor, în efectuarea unor quest-uri, competiție împotriva altor utilizatori, etc.

Grupurile de jucători pot avea caracter :

- temporar (echipe/teams/party)
- persistent (ghilde)

În general, orice grup are un conducător, care decide admiterea în sau excluderea din echipă a unui membru.

## **8) Termeni (acronime) des întâlniți în jocurile MMOG**

*AC* – armor class : valoarea numerică a „armurii” unui obiect

*Alt* – un caracter atașat contului unui utilizator altul decât caracterul principal

*AoE* – Area of Effect : termen asociat unei „vrăji” care afectează un număr mare de monștrii aflați într-o zonă largă

*AFK* – Away from keyboard. Jucătorul nu se mai află în fața calculatorului. De obicei sistemul de joc trece utilizatorul în aceasta stare după un timp de inactivitate al acestuia

*Aggro / Threat* – nivelul de amenințare pe care un jucător îl are față de un monstru din joc. De obicei într-un grup de mai mulți jucători , monstrul din joc va ataca jucătorul care are „Aggro”-ul cel mai ridicat. Exista mai multe modalități prin care acest nivel poate crește pentru un jucător, variind de la o simplă apropiere de monstrul respectiv până la atacuri directe asupra acestuia.

*Aggro Radius* – zona din jurul unui monstru în care dacă un jucător intră va genera „Aggro”

*Buff* – O vrajă benefică care este executată asupra unui caracter din joc

*Caster* – tip de jucător care pentru atacuri apelează de obicei la vrăji

*DOT* – Damage over time. Tip de damage care nu afectează cu întreaga valoare un caracter din momentul aplicării ci în decursul unui interval de timp.

*DPS* – Damage per second. Valoare a atacului instant al unei arme sau a unei vrăji pe secundă. De exemplu dacă o armă are o valoare asociată atacului și un timp de execuție , DPS-ul este raportul acestora.

*De-buff* – O vrajă cu efect negativ care este executată asupra unui caracter din joc

*Grinding* – rămânerea pentru un timp ridicat într-o anumită zonă din joc pentru a lupta cu același tip de monștri pentru a câștiga fie experiență , fie resurse care se obțin în urma înfrângerii acestora

*Guild* – grupare permanentă de jucători , de obicei cu un număr mare al acestora

*Health* - Sănătatea asociată caracterului unui jucător. Atunci când aceasta este epuizată , caracterul este incapacitat.

*Instanță* – o zonă din joc care se copiază pentru fiecare grup care intră în aceasta. Doar grupul care intră în „instanța” respectivă o poate accesa , aceasta nefiind vizibilă celorlalți utilizatori.

*Loot* – resursele care au fost obținute în urma înfrângerii unui monstru

*Mana* – atribut necesar pentru a executa vrăji sau atacuri care au asociate ca cerință acest atribut. Atunci când aceasta este epuizată nu se mai pot executa vrăjile sau atacurile care necesită o anumita cantitate din acest atribut.

*Newbie / Newb* – un termen care de obicei desemnează un jucător nou în joc. Câteodată este folosită și ca apelativ pentru un jucător care nu joaca foarte bine și face greșeli în joc.

*NPC* - non player-controlled character. Un caracter care este controlat ( parțial sau total) de către calculator.

*Party* – un grup de jucători care se asociază de regulă pentru îndeplinirea unor obiective comune.

*Pet* – un NPC care este controlat parțial de jucător

*PK* – Player killer. Un tip de jucător care de obicei are drept obiectiv înfrângerea celorlalți jucători prezenți în joc.

*PvE* - Player Vs. Environment. Acest termen descrie competiția între jucătorii din sistem și oponentii controlați de către acesta.

*PvP* – Player Vs. Player. Acest termen descrie competiția directă dintre jucătorii sistemului.

*Raid* – O grupare de jucători care are de obicei un număr mult mai mare decât un „Party”. Aceasta grupare se realizează de obicei atunci când este necesară îndeplinirea unui obiectiv foarte dificil care nu ar putea fi realizat de un număr mic de jucători.

*Tank* – tip de jucător melee care preferă atacurile directe cu arme

## **2.2 CATEGORII DE APLICAȚII MMO**

Din punct de vedere al jocului, există câteva categorii de tipuri de jocuri MMOG. Aceste categorii sunt următoarele :

- **MMORPG:** MMO Role-Playing Game – sunt cele mai faimoase și răspândite tipuri de jocuri MMOG. În acest tip de joc utilizatorul este transpus în rolul unui personaj virtual pe care trebuie să îl dezvolte pe parcursul jocului.
- **MMORTS:** MMO Real-Time Strategy – acest tip de joc combină tipurile de joc RTS cu o lume virtuală persistentă. De obicei în acest tip de joc , utilizatorul este pus în rolul unui conducător al cărui rol este să conducă o armată în batalii și să gestioneze resursele necesare pentru a-și conduce spre victorie armata
- **MMOFPS:** MMO First-Person Shooter – în acest tip de joc este combinat conceptul de lume virtuală cu cel de FPS accentul punându-se pe luptele directe la scară largă între jucători
- **MMOSG:** MMO Sports Game – în acest tip de joc , utilizatorii concurează în anumite sporturi cu mare răspândire cum ar fi fotbal , golf, baschet , fotbal american , etc
- **MMOMG:** MMO Management Game – acest tip de joc nu necesită în general prezența utilizatorului în mod permanent în lumea virtuală. În mod uzual utilizatorul se loghează de câteva ori pe săptămână în cont , stabilește ordinele și acțiunile care trebuiesc fie îndeplinite de către unitățile care le

controlează și strategia corespunzătoare pentru a câștiga. De obicei, acest tip de joc este jucat în interiorul unui browser

În continuare vor fi prezentate câteva exemple de jocuri MMOG care sunt foarte răspândite și au un număr ridicat de utilizatori.

### **World of Warcraft**

URL: <http://www.worldofwarcraft.com/>

Producător: Blizzard Entertainment

Data de lansare: lansarea pentru public a avut loc în 2004

3 expansiuni aparute ( The Burning Crusade , Wrath of The Lich King, Cataclysm)

Costuri: abonament lunar / gamecard ( 14.99 \$ / luna)

### **Descriere:**

World of Warcraft este poate cel mai cunoscut joc de tipul MMOG existent în prezent. Cu un număr de jucători activi ce depășește 10 milioane, se poate spune că este de asemenea și cel mai de succes MMOG din istoria acestui tip de joc. Jucătorul este introdus într-un univers fantasy și este pus să aleagă între două facțiuni aflate în conflict ( Alliance și Horde). Odată aleasă apartenența la una din aceste facțiuni jucătorul este introdus într-un univers fantasy bazat pe jocurile din seria « Warcraft » a cărui poveste i se va dezvălui pe măsură ce avansează în nivel. Jucătorul poate să lupte contra facțiunii adverse în diverse activități de tip PVP , să exploreze conținutul PVE variat care există în joc fie prin îndeplinirea de questuri sau prin parcurgerea numeroaselor instanțe aflate în joc sau să se ocupe cu câteva dintre numeroasele profesii existente. Jocul excelează prin dimensiunea și conținutul foarte bogat oferit atât din punct de vedere al lumii care poate fi explorată cât și al activităților pe care un jucător le poate îndeplini.



*Figura 2-1: World of Warcraft*

Facilități :

- Conținut PVE
  - Quest-uri
  - Instanțe (5-man , 10-man , 25-man și 40-man)
- Conținut PVP
  - 2 facțiuni opuse Alliance și Horde
  - Arena
  - Battlegrounds
- 10 Rase
  - Alliance : Draenei, Dwarves , Gnomes, Humans, Night Elves
  - Horde : Blood Elves, Orcs, Tauren , Trolls , Undead
- 10 Clase : Death Knight(Hero Class), Druid , Hunter , Mage , Paladin , Priest , Rogue , Shaman , Warlock , Warrior

- Sistem de profesii : Alchemy , Blacksmithing , Cooking , Enchanting , Engineering , First Aid , Fishing , Herbalism , Leatherworking , Mining , Skinning , Tailoring , Jewelcrafting , Inscription
- Sistem de health și mana
- Sistem de talente pentru fiecare clasă
- Sistem de achievements
- Sistem de transport public și mounts pentru jucători
- Posibilitate de alcătuire party de maxim 5 persoane
- Posibilitate de alcătuire raid de maxim 40 persoane
- Outdoor bosses
- Sistem de ghilde
- Action-House, Mail și Bank în-game
- 4 continente (Azeroth , Kalimdor , Northrend , Outland)
- Evenimente in-game
- Content nou adăugat periodic
- Forumuri

### **Lord of the Rings Online**

URL: <http://www.lotro.com/>

Producător: Turbine

Data de lansare: lansarea pentru public a avut loc în 2007

2 expansiuni apărute ( Mines of Moria, Siege of Mirkwood)

Costuri: free 2 play / microtranzacții

#### **Descriere:**

Lord of the Rings Online este construit pe universul creat de J.R.R. Tolkien în cărțile « The Hobbit » și « Lord of The Rings ». Jocul redă cu foarte mare fidelitate locațiile din Middle-Earth îndrăgite de fanii seriei și se bazează pe o grafică modernă și foarte detaliată , acest lucru contribuind foarte mult la atragerea jucătorilor. Acestora li se ofera posibilitatea urmării și implicării directe în povestea descrisă în cărțile « Lord of The Rings » prin îndeplinirea obiectivelor din lanțul de questuri « Epic » care este organizat sub forma de cărți. De asemenea jucătorii vor descoperi detalii despre lumea « Middle-earth »-ului și pot înfrunta multe dintre pericolele existente în universul « Lord of The Rings » prin îndeplinirea de questuri și instanțe. Jocul excelează prin grafica oferită și prin imersiunea ridicată a jucătorilor în universul virtual creat.



*Figura 2-2: Lord of The Rings Online*

Facilități :

- Conținut PVE
  - Quest-uri
  - Instanțe ( 6-man , 12-man , 24-man)
- Conținut PVP : deși nu exista facțiuni opuse , jocul ofera posibilitatea de PVP prin implementarea mecanismului de „Monster Play” unde fiecare jucător poate să aleagă să impersoneze un monstru având posibilitatea să joace contra monștrilor controlați de alți jucători într-o instanță specială
- 4 Rase : Man , Elf , Dwarf, Hobbit
- 9 Clase : Burglar, Captain, Champion, Guardian, Hunter , Lore-master , Ministrel , Rune-keeper, Warden
- Sistem de vocații : în funcție de vocația aleasă vor fi disponibile câte 3 profesii distincte :
  - Armorer : Prospector, Metalsmith, Tailor
  - Armsman : Prospector, Weaponsmith, Woodworker
  - Explorer : Forester, Prospector , Tailor
  - Historian : Farmer , Scholar , Weaponsmith



- Tinker : Prospector , Jeweller , Cook
- Woodsman : Forester , Farmer , Woodworker
- Yeoman : Farmer , Tailor , Cook
- Sistem de morală (health) și power (mana)
- Sistem de skills pentru fiecare clasă
- Sistem de individualizare a caracterului prin : deeds, traits , titles
- Sistem de muzică : odată cu nivelul 5, fiecare caracter poate să compună muzica proprie folosind instrumente specifice
- Sistem de transport public și mounts pentru jucători
- Posibilitate de alcătuire party de maxim 6 persoane ( fellowship)
- Posibilitate de alcătuire raid de maxim 24 persoane
- Action-House, Mail și Bank in-game
- Sistem de ghilde (kinship)
- Posibilitatea de explorare a Middle-Earth ( cu câteva locații bine cunoscute cum ar fi : Shire , Bree , Rivendell , Moria , Lorien)
- Evenimente in-game
- Conținut nou adăugat periodic
- Forumuri

### Warhammer Online

URL: <http://www.warhammeronline.com/>

Producător: Mythic Entertainment

Data de lansare: lansarea pentru public a avut loc în 2008

Costuri: abonament lunar / gamecard ( 14.99 \$ / luna)

### **Descriere:**

Warhammer Online este un joc în care accentul se pune în principal pe gameplay-ul de tip PVP. Jocul este unul de tip Realm vs Realm în care două facțiuni se luptă pentru a cuceri realm-urile inamicului. În acest tip de joc , utilizatorii fac parte din armate aparținând unei facțiuni și luptele între armate sunt date la scară largă cu participarea unui număr mare de jucători.

Jocul fiind axat în principal pe PVP, există un număr foarte mare de clase disponibile jucătorilor în funcție și de rasa aleasă, cu roluri bine definite : tank , melee și ranged dps , healer.



Figura 2-3: WarHammer Online

### Facilități :

- Conținut PVE
  - Public Quest
- Conținut PVP : Realm vs Realm prin :
  - Skirmish
  - Battlefields
  - Scenarii
- 2 facțiuni organizate sub forma a două armate :
  - Order
  - Destruction
- 6 Rase în funcție de facțiune :
  - Order : Dwarfs , Empire (Humans) , High Elves
  - Destruction : Greenskins ( Orcs , Goblins), Chaos (Humans), Dark Elves
- 24 de Clase (Cariere) în funcție de rasă ( câte 4 pentru fiecare rasa cu rolurile de Tank , Melee DPS, Ranged DPS și Healer) :
  - Dwarfs : Ironbreaker , Slayer , Engineer, Runepriest

- Empire : Knight of the Blazing Sun, Witch Hunter , Bright Wizard , Warrior Priest
  - High Elves : Swordmaster , White lion , Shadow Warrior , Archmage
  - Greenskins : Black Orc, Choppa , Goblin Squig Herder , Goblin Shaman
  - Chaos : Chosen , Marauder , Magus , Zealot
  - Dark Elves : Black Guard , Witch Elf , Sorceress/Sorcerer, Disciple of Khaine
- Sistem de health și action points (mana)
  - Collision detection avansat între personajele jucătorilor creând astfel noi mecanici de joc
  - Sistem de skills pentru fiecare clasă bazat pe specializări
  - Sistem de transport public și mounts pentru jucători
  - Posibilitate de alcătuire party de maxim 6 persoane
  - Posibilitate de alcătuire raid (warbands) de maxim 24 persoane
  - Action-House, Mail și Bank în-game
  - Conținut nou adăugat periodic
  - Forumuri

## **2.3 TEHNOLOGII CLIENT-SIDE**

Așa cum s-a precizat și anterior, un aspect foarte important, atunci când vorbim de simulările 3D, este acela al imersiunii utilizatorului în cadrul lumii virtuale. Pentru a realiza acest lucru, trebuie să existe pe partea de client o redare grafică care folosește tehnologii 3D moderne pentru a avea :

- o reprezentare foarte detaliată a lumii virtuale, acest lucru presupunând reprezentarea obiectelor din spațiul virtual cu un număr ridicat de poligoane;
- efecte vizuale spectaculoase cum ar fi :
  - fum
  - explozii
  - umbre dinamice
  - texturi detaliate, etc

OpenGL și Direct3D, principalele biblioteci dedicate dezvoltării de aplicații grafice 3D, permit crearea și redarea de spații 3D în care un observator virtual poate naviga folosind mouse-ul sau tastatura. Scena este descrisă prin modelele 3D ale obiectelor, poziționarea acestora și a observatorului în scena, proprietățile de material ale suprafețelor obiectelor, poziționarea și proprietățile surselor de lumina, texturile care trebuie aplicate pe suprafețele obiectelor, și altele.

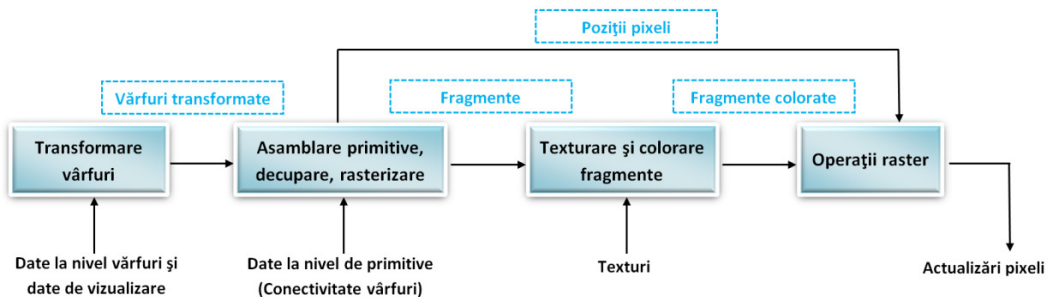
### 2.3.1 BANDA GRAFICĂ (GRAPHICS PIPELINE)

Componenta grafică a oricărei lumi virtuale este probabil domeniul cel mai dinamic, aflat mereu în căutare de mai multă performanță, atât la nivel de arhitecturi hardware și software cât și la nivel de algoritmi [PET97a] și detalii de implementare.

Modelul scenei 3D, descris de o aplicație apelând funcțiile interfeței de programare (API) OpenGL / Direct3D, poate fi compus din simple primitive grafice (linie, triunghi, poligon), primitive grafice conectate (set de triunghiuri conectate, set de poligoane conectate), obiecte 3D, cum ar fi sfera, și suprafețe de forma liberă (suprafețe Bezier, B-spline, Nurbs). Bibliotecile OpenGL și Direct3D discretizează suprafețele obiectelor 3D și pe cele de forma liberă aproximându-le printr-o rețea spațială conectată de fațete poligonale (patrulaterale sau triunghiuri).

#### 2.3.1.1 Banda grafică fixă

Transformarea descrierii scenei 3D într-o imagine afișată pe ecran este realizată automat de API-urile 3D, printr-un lanț de prelucrări fixe. Lanțul de prelucrări constă din etapele redate în figura următoare:



*Figura 2-4: Banda grafică fixă (clasică)*

Deoarece fiecare etapă primește intrări de la etapele anterioare și transmite rezultate pentru etapele următoare, secvența de operații este asemănată cu o “bandă de asamblare” într-un proces de fabricație. Deoarece obiectul fabricat în acest caz este imaginea unei scene 3D, banda este denumită „banda grafică” („graphics pipeline”).

Operațiile din **banda grafică fixă**, redată în figura de mai sus, sunt aceleași pentru toate primitivele unei scene 3D.

În cele ce urmează vom explica pe scurt aceste operații.

#### *2.3.1.1.1 Transformarea vârfurilor*

Operațiile din această etapă se efectuează asupra fiecărui vârf. Datele de ieșire ale etapei sunt folosite în etapa următoare a benzii grafice. Astfel, poziția vârfului, care este specificată în spațiul obiect, este transformată în spațiul de proiecție (spațiul de decupare) pe baza matricei de modelare-vizualizare-proiecție (specificată ca un parametru global de vizualizare). Dacă poziția este singura dată asociată vârfului, atunci toți pixelii primitivei din care face parte vor fi afișați în aceeași culoare, specificată ca o dată globală de vizualizare (iluminare constantă).

#### *2.3.1.1.2 Asamblare primitive, decupare, rasterizare*

În această etapă vârfurile sunt conectate în primitive și primitivele sunt supuse operației de decupare la marginile volumului vizual canonic (view frustum clipping).

Din decupare rezultă noi vârfuri. Orice nou vârf este situat pe o latură a primitivei decupate. În funcție de poziția noului vârf pe latură decupată, se calculează prin interpolare datele asociate noului vârf: culoarea, coordonatele textură și coordonata ceții.

Pozițiile vârfurilor primitivelor rezultate după decupare sunt transformate din spațiul de decupare în spațiul ecran apoi primitivele sunt “rasterizate”, adică descompuse în rastru de puncte corespunzător spațiului discret al suprafeței de afișare. Rezultatul rasterizării este un set de adrese de pixeli și un set de “fragmente” de primitivă care trebuie să fie afișate în pixeli.

În acest punct, trebuie făcută foarte clar distincție între termenul de “pixel” și “fragment”.

Un pixel reprezintă conținutul framebuffer-ului la o locație specifică cum ar fi culoarea, adâncimea și alte valori asociate cu locația respectivă.

Un fragment este un “potential pixel”. Dacă fragmentul trece testele de rasterizare (în etapa de Operații raster) el va actualiza pixelul caruia îi este asociat în framebuffer.

#### *2.3.1.1.3 Texturare și colorare fragmente*

Fiecare fragment are asociată o adresă de pixel, o valoare de adâncime (corespunzătoare coordonatei Z a fragmentului, ce reprezintă distanța sa față de observatorul virtual al scenei 3D) și o culoare. De asemenea, fragmentului îi mai pot fi asociate unul sau mai multe seturi de coordonate textură și coordonata de ceață. Coordonatele textură și coordonata ceații se obțin prin interpolarea valorilor atașate vârfurilor primitivei din care face parte fragmentul.

#### **Modelul de colorare (shading)**

OpenGL și Direct3D permit aplicațiilor specificarea modului de calcul al culorii fragmentelor: aceeași culoare pentru toate fragmentele unei primitive, preluată de la unul dintre vârfuri, sau calculul culorilor fragmentelor prin interpolarea culorilor asociate vârfurilor după etapa de transformare a vârfurilor (modelul Gouraud).

#### **Texturare**

Coordonatele textură permit calculul unei culori prin accesarea unei texturi (1, 2 sau 3-dimensionale). Culoarea textură se combină cu culoarea fragmentului folosind diverși operatori. Transparența texturii poate fi controlată prin canalul alfa.

În cazul multi-texturării (fragmentul are asociate mai multe seturi de coordonate textură), culoarea fragmentului se combină succesiv cu culorile extrase din diferite texturi. Multi-texturarea permite obținerea de efecte care măresc realismul scenei sintetizate.

#### *2.3.1.1.4 Operații raster*

Operațiile raster se efectuează asupra fiecărui fragment înainte de actualizarea buffer-ului imagine (framebuffer). Ele pot fi specificate prin interfețele de programare OpenGL și Direct3D. În această etapa sunt eliminate părțile nevizibile ale primitivelor care alcătuiesc scena 3D, pe baza testului de adâncime (algoritmul z-buffer). Alte operații raster sunt: testul de apartenență a pixelului la contextul OpenGL curent, testul de decupare la marginile unei ferestre din zona de afișare (scissor test), testul alfa, testul stencil, testul de combinare (blending) precum și diverse operații logice între culoarea fragmentului și culoarea existentă în buffer-ul imagine.

Culorile finale ale fragmentelor care “trec” de operațiile raster sunt afișate în pixelii corespunzători.

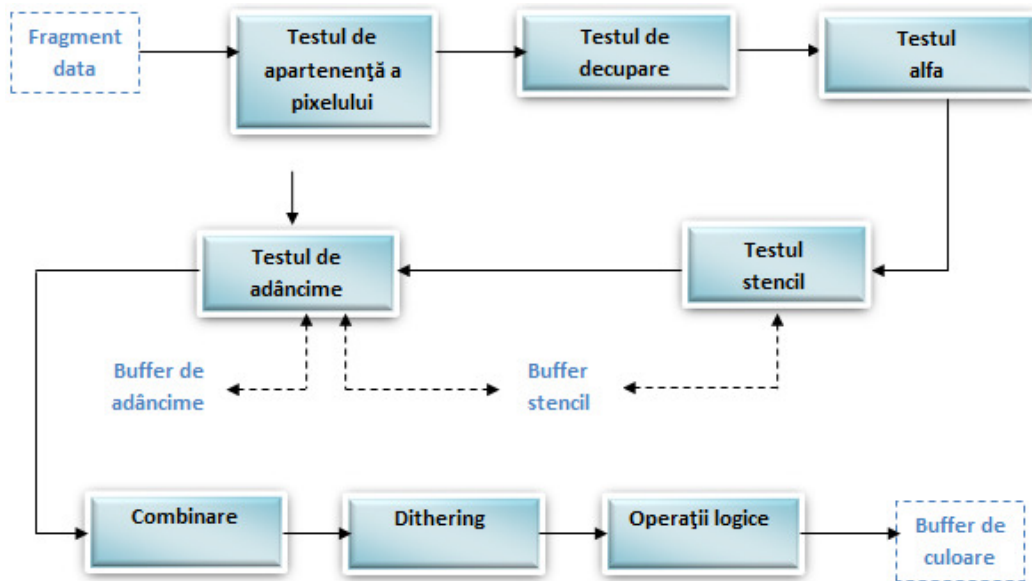


Figura 2-5: Operații raster

### 2.3.1.2 Banda grafică programabilă

Modul în care sunt implementate operațiile din banda grafică fixă depinde de suportul hardware oferit de unitatea (placa) grafică a calculatorului. Unitatea de Prelucrare Grafică poate implementa diferite nivele de accelerare grafică, de la un sistem raster capabil să rasterizeze linii și poligoane, până la un sistem complex de procesoare capabile să transforme și să prelucreze date geometrice, cu un înalt grad de paralelism. Astfel, implementarea fiecărei comenzi OpenGL/Direct3D este împărțită în mod optim între CPU (Unitatea de Prelucrare Centrală a Calculatorului) și GPU, în funcție de capacitatea de prelucrare a GPU-ului.

Ultimele versiuni ale bibliotecilor OpenGL/Direct3D includ suportul necesar pentru implementarea conceptului de **bandă grafică programabilă**. Programatorul are posibilitatea de a furniza programele care sunt executate la nivel de vârf, respectiv fragment, înlocuind prelucrările standard din banda grafică fixă. Aceste programe se numesc **vertex shader** și **pixel shader** deoarece principala lor utilizare este în calculul culorii de afișare a pixelilor. Ele pot fi implementate în limbaje de nivel înalt: **Cg (C for Graphics)** [FEK06] oferit de NVIDIA și **HLSL (High Level Shading Language)** dezvoltat de Microsoft pentru a fi utilizat împreună cu

Direct3D (de fapt, Cg și HLSL reprezintă un același limbaj, rezultat din colaborarea dintre NVIDIA și Microsoft), analog cu **GLSL (OpenGL Shading Language)**.

Prin intermediul programelor vârf (vertex shader) și pixel (pixel shader) programatorul poate implementa modalități diverse de transformare a vârfurilor, diferite de cele standard, de exemplu pentru realizarea animației, precum și calcule de iluminare mai exacte la nivel de fragment. De exemplu, în loc să se interpoleze culorile vârfurilor pentru obținerea culorilor fragmentelor interioare, se poate calcula o culoare pentru fiecare fragment folosind normala interpolată (modelul Phong), sau o normală citită dintr-o „hartă de normale”, etc.

Primul model de programare a Unității de Prelucrare Grafice bazat pe vertex shader și pixel shader (“shader model 1”) a apărut odată cu generația a 4-a de GPU (începând cu 2003: NVIDIA GeForce FX, ATI Radeon 9700 ). Generația anterioară suporta programabilitatea numai la nivel de vârf. Penultimul model de programare a Unității de Prelucrare Grafică, «**shader model 4**» conține încă o etapă programabilă în banda grafică, implementată prin “**geometry shader**”, care permite crearea de noi primitive în banda grafică. Acest model a fost introdus în noiembrie 2006 de Direct3D10 și GLSL1.2 cu suport hardware GPU NVIDIA 8800GTX (2006) și Radeon 2900XT (2007). Ultimul model «**shader model 5**», apărut în 2008, conține etape programabile dedicate tehnicii de „tessellation” și pentru calcule GPGPU.

Unitățile de Prelucrare Grafică moderne, cu care sunt dotate tot mai multe calculatoare personale în prezent, pot executa integral operațiile din banda grafică. Execuția unei aplicații OpenGL/Direct3D pe un astfel de calculator este mult mai rapidă, redarea scenelor 3D cu texturi sofisticate, lumini și umbre putând fi realizată în timp real. Aceasta deoarece CPU-ul este eliberat de operațiile necesare la nivelul fiecărui vârf și pixel, fiind disponibil pentru alte operații ale aplicației. Totodată, execuția de către GPU a operațiilor din banda grafică este mult mai rapidă, atât datorită specializării procesoarelor pentru operațiile respective cât și datorită paralelizării ridicate a prelucrărilor. Datorită specializării hardware și paralelismului operațiilor la nivelul GPU, sinteza imaginilor foto-realiste în timp real a devenit posibilă. De asemenea, animația personajelor și alte metode de animație pot beneficia din plin de puterea de calcul a GPU-ului.



Figura următoare reda schematic etapele de prelucrare din “banda grafică programabilă”, conform modelului de programare a Unității de Prelucrare Grafică 4 (“shader model 4”) [FAT08].

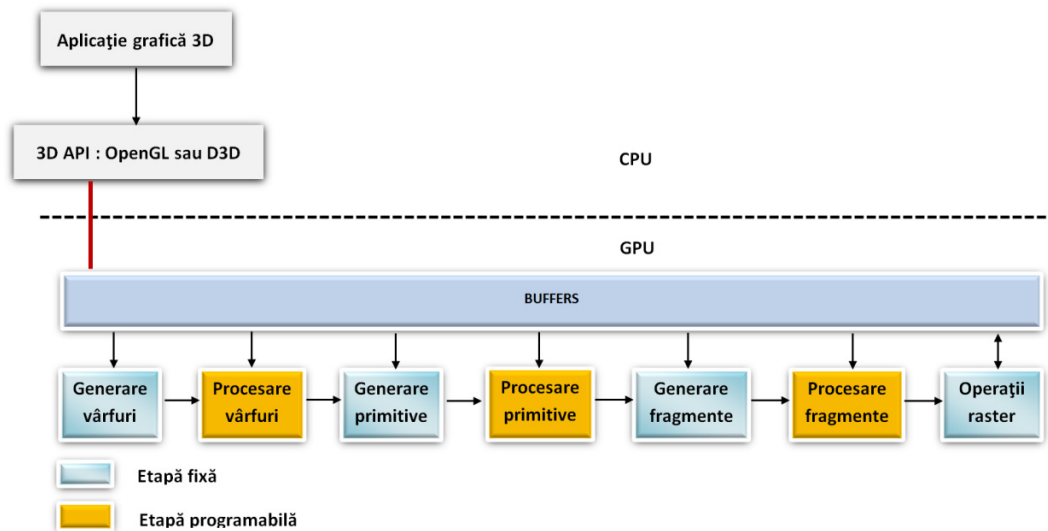


Figura 2-6: Banda grafică programabilă

Banda grafică programabilă conține 7 etape, dintre care 4 sunt fixe și 3 sunt programabile.

#### 2.3.1.2.1 Generare vârfuri

Fiecare primitivă grafică este definită într-o aplicație OpenGL/Direct3D printr-un set de vârfuri. Pornind de la aceste vârfuri, în această etapă, se creează un șir (stream) de înregistrări “vârf” care se transmit mai departe în banda grafică. Fiecare înregistrare conține coordonatele unui vârf, împreună cu alte date asociate vârfului (culoare, normală, coordonate textură). Acestea sunt intrările pentru “vertex shader”, care este executat în paralel, pentru diferite înregistrări vârf, de mai multe procesoare ale GPU.

#### 2.3.1.2.2 Prelucrare vârfuri

Această etapă este programabilă și poate fi implementată printr-un program vârf (“vertex shader”). Atunci când este activ un program vârf, el este apelat pentru fiecare vârf produs în etapa anterioară. Principala operație realizată într-un program vârf este transformarea coordonatelor vârfului în spațiul de proiecție și calculul

culorii vârfului. Aici însă pot fi programate și alte transformări asupra vârfului, de exemplu pentru animația personajelor (“skinning”), “morphing”, simularea valurilor, etc. În afara datelor primite din banda grafică, programul vârful utilizează și date din buffer-ele globale ale memoriei GPU, de exemplu, poziția observatorului, matricea de proiecție, poziția sursei de lumină și alte date specifice aplicației.

Ieșirea programului vârful este o înregistrare “vârful de ieșire”, conținând coordonatele transformate ale vârfului, culoarea sa, una sau mai multe coordonate textură.

La un moment dat, există un singur program vârful, care se execută asupra tuturor vârfulor primitivelor.

#### *2.3.1.2.3 Generare primitive*

Această etapă nu este programabilă. Scopul său este de a asambla vârfulurile în primitive, folosind topologia fiecărei primitive definită în aplicație. Este generat un șir ordonat de înregistrări “primitivă”, rezultate prin concatenarea înregistrărilor vârful produse de programul vârful.

#### *2.3.1.2.4 Prelucrare primitive*

Este o etapă programabilă, ce poate fi implementată printr-un program “primitivă” (“primitive shader” sau “geometry shader”). Acest program acționează asupra unei singure primitive (primitivă de intrare), care poate fi: un singur vârful, o linie (2 vârfuluri) sau un triunghi. De asemenea, datele de intrare pot fi vârfulurile primitivelor adiacente cu primitiva de intrare: două vârfuluri pentru o linie (extremitățile segmentelor adiacente), trei vârfuluri pentru un triunghi (vârfulurile suplimentare ale triunghiurilor adiacente). Programul poate produce un nou șir ordonat de primitive: șir de vârfuluri (PointStream), șir de linii (LineStream) sau șir de triunghiuri (TriangleStream).

Atunci când este activ un program “primitivă”, el este apelat o dată pentru fiecare primitivă generată în etapa anterioară.

Exemple de algoritmi care pot fi implementați în programul primitivă sunt: sisteme de particule dinamice, generare de “sprite”-uri, generarea volumului de umbră al primitivei, generarea terenurilor, s.a.

#### *2.3.1.2.5 Generare fragmente*

În această etapă este efectuată decuparea primitivelor la marginile volumului vizual și apoi rasterizarea lor, rezultatul fiind un șir de înregistrări “fragment”, care sunt

transmise procesoarelor de prelucrare a fragmentelor. Fiecare înregistrare fragment conține poziția pixelului în care va fi afișat fragmentul, coordonata z a fragmentului și valorile rezultate din interpolarea datelor atașate vârfurilor (culori, coordonate textură).

#### *2.3.1.2.6 Prelucrare fragmente*

Prelucrarea din această etapă poate fi programată printr-un program “fragment” și are drept scop calculul culorii și al opacității unui fragment. În acest scop pot fi folosite texturi 1D, 2D sau 3D, se pot simula efecte de reflexie și refracție a luminii de către suprafețele obiectelor din scena 3D.

#### *2.3.1.2.7 Operații raster*

În această etapă sunt eliminate fragmentele obturate de suprafețe aflate mai aproape de observator (algoritmul z-buffer), se combină culoarea fragmentului și culoarea curentă a pixelului folosind opacitatea și diverși operatori raster.

Programele “shader” folosesc, nu numai datele primite prin șirurile de înregistrări ale benzii grafice, ci și date memorate în buffere globale existente în memoria GPU (memoria video). Aceste buffere sunt inițializate cu datele necesare programelor “shader”, furnizate de programul de aplicație.

### **2.3.2 ORGANIZAREA ȘI REPREZENTAREA SCENEI 3D**

Alegerea unei structuri de date adecvate pentru reprezentarea scenei este extrem de importantă pentru optimizarea operațiilor necesare redării sale. Astfel, una dintre optimizările cele mai importante este excluderea din banda grafică a obiectelor care nu sunt vizibile în imagine. Într-o aplicație OpenGL această excludere este realizată automat la nivelul fiecărei primitive geometrice, în banda grafică, prin operația de decupare („culling”) la marginile volumului vizual definit de programator. O metodă mai eficientă ar fi excluderea la nivel de obiect și nu de primitive în care a fost descompus obiectul. Dar și mai eficient este să se poată exclude un întreg grup de obiecte nevizibile. Evident, între obiectele grupului trebuie să existe o relație spațială sau logică. Ideea de a grupa obiectele este utilă și în privința transformărilor. De multe ori, aceeași transformare se aplică tuturor primitivelor în care este descompus obiectul atunci când este transmis benzii de redare.

Dezideratele menționate mai sus au condus la ideea managementului scenei 3D la un nivel superior celui oferit de OpenGL sau Direct3D.

### **2.3.2.1 Managementul ierarhic al scenei**

Structura ierarhică utilizată de regulă pentru reprezentarea scenei este de tip **arbore**, în care fiecare nod are un părinte și orice număr de copii, iar fiecare copil are un singur părinte. O structura mai generală este de tip **graf aciclic orientat (DAG)**, care permite reprezentarea partajării obiectelor (un nod poate avea mai mulți părinți). Structura ierarhică folosită pentru reprezentarea scenei este cunoscută sub numele de „**graful scenei**” (**scene-graph**).

**Nodurile frunză** conțin primitivele care sunt transmise benzii grafice dar și alte date care au un efect final asupra scenei redată în imagine, de exemplu lumini sau sunet. Celelalte noduri au scop de grupare, de aceea sunt denumite **noduri grup**.

Nodurile arborelui conțin **Informații spațiale** și **informații semantice** ca: transformări, volume încadratoare, Informații de redare și de animație. Transformările sunt utilizate pentru poziționarea, orientarea și scalarea obiectelor din ierarhie. Volumele încadratoare sunt utilizate pentru testul de vizibilitate (apartență la volumul vizual) și în detecția coliziunilor (testele de intersecție). Informațiile de redare atașate unui nod pot fi, de exemplu, o sursă de lumină care luminează obiectele din subarborele nodului, nivelul de ceață, texturile folosite la redare și altele.

Exemple de tipuri de noduri în graficul unei scene sunt: "cub", "sfera", "text", "teren", "skyBox", "arbore octal", "mesh", "sursă de lumină", "transformare", "cameră", "sistem de particule", etc. Indiferent de tipul sau, fiecărui nod îi este asociată o poziție în spațiul scenei.

Modificările care apar în scena 3D pe parcursul execuției se traduc în modificări ale atributelor nodurilor arborelui. O schimbare a unui nod afectează subarborele a cărui rădăcină este. Modificarea trebuie transmisă tuturor nodurilor subarborelui. Astfel, schimbarea poziției unui obiect reprezentat printr-un nod intern are ca efect schimbarea transformărilor care afectează nodurile descendente în graficul scenei, precum și schimbarea volumului încadrator al fiecărui nod.

### **2.3.2.2 Reprezentări bazate pe sortare spațială**

Reprezentarea ierarhică a scenei permite evitarea transmiterii în banda grafică a unui număr mare de obiecte nevizibile. Totuși, un număr mare de primitive care sunt incluse în volumul vizual dar se suprapun în imagine sunt rasterizate, iar

problema vizibilității este rezolvată la nivel de pixel, prin testul de adâncime (coordonatele z, în algoritmul z-buffer).

Ideea de bază a sortării spațiale a obiectelor scenei este de a împiedica, pe cât posibil, “scrierea” de mai multe ori a unui pixel (în buffer-ul imagine) în procesul de formare a imaginii unei scene 3D. Termenul “depth complexity” se referă la numărul de scrieri al unui pixel în acest proces. Se dorește ca aceasta complexitate să fie 1. Cu cât este mai mare cu atât timpul necesar formării imaginii este mai mare.

#### *2.3.2.2.1 Arborele BSP (Binary Space Partitioning)*

Arborele BSP este o metodă frecvent folosită pentru sortarea spațială a obiectelor scenei 3D. Arborele BSP, prin care se reprezintă scena 3D, se obține prin partiționarea binară, recursivă, a scenei, folosind plane de partiționare.

Reprezentarea este independentă de poziția observatorului, de aceea este convenabilă în aplicațiile în care scena nu se modifică pe parcursul execuției aplicației, arborele fiind construit într-o etapă de pre-procesare.

### **Construirea arborelui**

De regulă, planele de partiționare sunt planele poligoanelor în care a fost descompusă scena 3D. Construirea scenei poate începe cu oricare dintre poligoane sau cu unul ales astfel încât numărul de divizări de poligoane la construirea arborelui să fie redus. Nodul rădăcină reprezintă întreaga scenă. Lui  $i$  se atașează poligonul de start. Față de planul acestui poligon (planul de partiționare), celelalte poligoane sunt clasificate în două liste: lista poligoanelor „din față” și lista poligoanelor „din spate”. Clasificarea se face pe baza orientării normalei la planul poligonului (poligoanele sunt contururi orientate). Astfel, poligoanele „din față” se află în semispațiul din fața planului de partiționare (semispațiul spre care este orientată normala la planul poligonului), iar poligoanele „din spate” sunt conținute în semispațiul spate. Poligoanele intersectate de planul de partiționare sunt divizate, rezultând alte poligoane care se introduc în lista poligoanelor „din față” sau „din spate” în funcție de poziționarea lor față de planul de partiționare. În arborele BSP, poligoanele „din față” se vor afla în subarborele „față” (sau „stânga”) al nodului rădăcină iar cele „din spate” în subarborele „spate” (sau „dreapta”) al nodului. Funcția de construire a arborelui se apelează recursiv pentru cele două liste de poligoane, ieșirea din recursivitate având loc atunci când lista de poligoane de intrare este vidă. Fiecărui nod al arborelui îi este atașat un poligon al scenei.

### **Afișarea „din spate în față”**

Afișarea scenei reprezentate prin arbore BSP este, de regula, de tip „din spate în față”, ținând cont de poziția observatorului. Funcția de afișare este apelată inițial pentru nodul rădăcină apoi, recursiv, pentru toate celelalte noduri, până la nodurile frunză. Astfel, dacă observatorul se află în semispațiul „spate” al planului (poligonului) nodului curent și planul intersectează volumul vizual, atunci se vor afișa mai întâi poligoanele din subarborele „față”, apoi poligonul din nod și la sfârșit poligoanele din subarborele „spate”. Dacă planul nodului curent nu intersectează volumul vizual, atunci subarborele „față” al nodului nu intersectează volumul vizual (este invizibil) și nu este procesat, fiind eliminat complet din banda grafică, la fel ca poligonul din nodul curent. Dacă observatorul se află în semispațiul „față” al nodului curent, ordinea de afișare este inversă.

### **Afișarea „din față în spate”**

Afișarea „din spate în față” asigură redarea corectă a scenei 3D dar numărul de suprascriseri ale pixelilor în procesul de redare poate fi semnificativ. Din această cauză algoritmul poate să nu fie suficient de rapid pentru redarea în timp real.

Arborele BSP permite afișarea scenei „din față în spate”, adică începând cu poligoanele cele mai apropiate de observator. Odată scris un pixel, el nu mai trebuie rescris în procesul de rasterizare a altor primitive, datorită corectitudinii sortării. În acest scop se folosesc diferite metode care folosesc măști la nivel de pixel.

## **2.3.3 TEXTURARE**

Utilizarea modelelor de iluminare pentru redarea suprafețelor obiectelor 3D nu este suficientă pentru obținerea de imagini realiste. Suprafețele reale nu sunt perfect netede și lucioase, efect produs prin iluminare și colorare. Ele au adesea detalii care, folosind numai iluminarea, pot fi redate numai prin rafinarea modelului lor geometric (creșterea numărului de poligoane din reprezentarea lor). Astfel de detalii pot fi: neregularități (rugozități), pete, zgârieturi, etc. Problema redării detaliilor este rezolvată prin utilizarea texturilor, nefiind necesară creșterea complexității modelului geometric.

În utilizarea normală, termenul de “textură” se referă la calitățile suprafeței unui obiect, din punctul de vedere al rugozității acesteia. În sinteza imaginilor pe

calculator însă, o textură reprezintă o colecție de culori pentru pixeli, memorată în general sub forma unei imagini, utilizată pentru a simula din punct de vedere vizual textura unui obiect.

Texturile utilizate sunt, în general, imagini bidimensionale care codifică diverse informații despre suprafața corespunzătoare. În cazul cel mai simplu, acestea sunt utilizate pentru a memora culoarea suprafeței în diverse puncte. De exemplu, o imagine care redă textura lemnului poate fi utilizată pentru a simula suprafețele obiectelor cu textură lemnoasă. În același mod pot fi simulate extrem de multe obiecte, cum ar fi, straturile de rocă din stânci, marmura de pe podele, etc.

Prin utilizarea mai multor texturi în paralel pentru redarea aceleiași suprafețe, fiecare în alt scop, se pot obține efecte mult mai complexe. De exemplu, este posibilă modificarea coeficientului de reflexie a luminii pe suprafața unui obiect care este parțial mat, parțial reflectant. Sau, este posibilă alterarea normalelor la suprafață pentru a simula obiectele cu suprafețe rugoase fără a fi necesară adaugarea detaliilor prin poligoane suplimentare.

Cele mai multe texturi sunt reprezentate de fapt de matrici care conțin în celule culori, codificate RGB (Red, Green, Blue). Fiecare astfel de valoare poartă numele de element al texturii, sau **texel**. Texelii sunt adresati într-un spațiu parametric bi-dimensional,  $(u,v)$ , cu  $0 \leq u,v \leq 1$ , numit **spațiul textura**.

Atunci când este necesară aplicarea unei texturi pe un obiect din spațiul 3D, adresele texelilor trebuie transformate în spațiul obiectului, după care pot fi transformate în coordonatele ecranului. În mod uzual însă, această mapare se efectuează invers, adică pentru fiecare pixel de pe ecran este calculată o coordonata textura,  $(u,v)$ , care se folosește la calculul unei culori din matricea textură. Metoda se numeste **mapare inversă**. Deoarece o pereche  $(u,v)$  nu poate selecta în mod unic un element al matricei textura, culoarea extrasă din textură se calculează printr-o **metoda de filtrare a texturii**. Filtrarea este în general realizata printr-o mediere biliniară între texeli vecini.

Maparea inversă presupune stabilirea unei corespondențe între punctele suprafeței 3D texturate și spațiul parametric al texturii. Această corespondență este definită matematic numai pentru anumite tipuri de suprafețe, care au o reprezentare parametrică: sfera, cilindrul, planul, suprafețele bicubice. Funcțiile matematice folosite în acest scop se numesc **funcții de mapare standard**. Orice suprafață 3D este descompusa de sistemul grafic, în vederea redării sale, într-o plasă poligonală. Astfel, texturarea se aplică primitivelor geometrice în care a fost descompusă

suprafața, triunghiuri sau poligoane, la momentul rasterizării lor. Problema mapării inverse apare deci pentru fragmentele acestor primitive.

Aplicarea texturilor pe suprafețele 3D este aproximată astfel:

- se calculează o adresă textură,  $(u,v)$ , pentru fiecare vârf al rețelei poligonale care aproximează suprafața; pentru aceasta se aplică **maparea în doi pași** ;
- pentru fiecare fragment al unui poligon, se calculează o adresa textură,  $(u_f, v_f)$ , prin interpolarea adreselor textură asociate vârfurilor sale. Pentru rezultate corecte este necesară efectuarea unei interpolări perspective. Interpolarea liniară este mai rapidă dar poate produce defecte ;
- în cazul utilizării mai multor texturi la redarea unui poligon, fiecărui vârf îi sunt asociate un același număr de perechi de coordonate textură  $(u_i, v_i)$ ,  $1 \leq i \leq n$ , pentru accesarea celor  $n$  texturi;  $n$  este în general limitat la 8.

### **Utilizarea coordonatelor textură**

Cu toate că spațiul parametric al texturii este adresat prin valori  $(u,v)$  cuprinse între 0 și 1, valorile  $(u,v)$  asociate vârfurilor pot fi în afara acestui domeniu. În procesul de rasterizare ele sunt interpolate pe suprafața primitivei, rezultând coordonate  $(u,v)$  asociate fragmentelor. Acestea trebuie transformate în valori  $0 \leq u', v' \leq 1$ .

Cele două moduri standard de transformare a coordonatelor textură în domeniul  $[0,1]$  sunt denumite “clamping” și “wrapping (sau “repeating”). Ele permit utilizarea eficientă a texturilor și obținerea de efecte interesante.

### **Filtrarea texturii**

O textură este definită pe o latică discretă de puncte în timp ce coordonatele textură sunt valori reale. În general, o pereche  $(u,v)$  nu corespunde unui texel al texturii, ci selectează o pereche  $(i', j')$ , unde  $i'$  și  $j'$  sunt numere reale. Metoda de calcul a culorii dintr-o textură, corespunzătoare unei perechi  $(u,v)$ , se numește “filtrarea texturii”.

Sunt două moduri standard de filtrare:

- Se consideră culoarea texelului cel mai apropiat în spațiul laticii (texelul cu indicii  $(i, j) = ([i' + 0.5], [j' + 0.5])$ , unde  $[x]$  reprezintă cel mai mare întreg mai mic sau egal cu  $x$ ).
- Filtrarea liniară : se efectuează o interpolare biliniară între culorile celor mai apropiați texeli în spațiul laticii.

Fiecare dintre cele două metode de filtrare are avantaje și dezavantaje. Primul tip de filtru este mai rapid dar poate produce defecte în imagine, în funcție de textură.



Dacă se preferă calitatea, se va alege filtrarea liniară, pe seama scăderii vitezei de generare a imaginii. Un efect secundar al filtrării liniare este de încetșoare a liniilor.

### **Transparența și opacitate**

Imaginea textură poate avea o componentă adițională, numită “canalul alfa”, prin care se controlează transparența sau opacitatea texturii aplicate. Fiecare texel este în acest caz reprezentat prin 4 componente: (R,G,B,A), unde A reprezintă opacitatea. Dacă  $A = 1$ , atunci texelul este complet opac. Dacă  $A=0$ , texelul este complet transparent. Pentru valori ale lui A cuprinse între zero și unu, culoarea texelului modulează culoarea fragmentului, determinată pe baza geometriei.

Interfețele de programare OpenGL și Direct3D permit specificarea modului de combinare între culoarea unui fragment, determinată pe baza geometriei primitivei, și culoarea rezultată din filtrarea texturii (texturilor), folosind operatori care se aplică culorilor (RGB) și operatori care se aplică opacităților (A).

#### **2.3.3.1 Environment mapping (maparea mediului înconjurător pe suprafețe)**

Redarea mediului înconjurător pe suprafețele obiectelor se poate aproxima printr-o metodă specială de aplicare a texturilor, cunoscută sub numele de “Environment mapping”. Mediul înconjurător este codificat într-o imagine, numită “harta mediului înconjurător”(environment map). Aceasta imagine se aplică sub forma unei texturi, pe suprafața unei sfere sau a unui cub, în centrul căreia / căruia se află obiectul. Dacă obiectul este reflectant (suprafața sa reflectă specular lumina incidentă) atunci harta textură este reflectată de suprafața obiectului. Dacă obiectul este transparent, atunci harta textură este redată pe suprafața obiectului prin fenomenul de refracție. În ambele cazuri, imaginea redată pe suprafața obiectului este dependentă de poziția observatorului, raza reflectată/refractată fiind calculată în funcție de poziția observatorului și normala la suprafață.

#### **2.3.3.2 Simularea detaliilor fizice ale suprafețelor**

Dupa cum am arătat mai înainte, calculul reflexiei luminii într-un punct al unei suprafețe utilizează normala la suprafața.

În cazul unei suprafețe netede, normala nu variază brusc în punctele sale vecine. Să ne imaginăm însă un perete alcatuit din cărămizi, sau suprafața unei stânci. În astfel de cazuri, exista variații bruște ale normalei în puncte vecine ale suprafeței. Pentru redarea exactă a unor astfel de suprafețe, folosind numai modelul de iluminare, este

necesar ca suprafața sa fie modelată geometric prin foarte multe poligoane. Astfel, fiecare cărămidă și fiecare spațiu dintre cărămizi ar trebui modelate prin poligoane. Modelul geometric devine și mai complicat dacă vrem să redăm o suprafață rugoasă, cum ar fi aceea a unei pietre.

### 2.3.3.2.1 *Bump mapping*

În loc de a modela geometric detaliile fizice ale unei suprafețe, putem să le simulăm, înlocuind, în calculul luminii reflectate, normala la suprafață fără detalii (de exemplu, peretele neted, fără cărămizi) cu o normală care ar corespunde suprafeței cu detalii (peretele cu cărămizi). Ideea a aparținut lui Blinn ([BLI78]), care a și introdus metoda cunoscută sub numele de “Bump mapping”. El a modelat matematic suprafața cu detalii fizice folosind o funcție de perturbare definită în fiecare punct al suprafeței fără detalii.

Astfel, dacă suprafața este definită parametric printr-o ecuație  $Q(u,v)$ , atunci normala la suprafața într-un punct  $(u,v)$  este dată de produsul vectorial dintre vectorii tangenți la suprafață în punctul  $(u,v)$ :

$$N = Q^u \times Q^v$$

unde  $Q^u$  și  $Q^v$  sunt derivatele parțiale în raport cu  $u$  și  $v$ . Suprafața perturbată este definită astfel:

$$Q'(u, w) = Q(u, w) + P(u, w) \frac{\bar{N}(u, w)}{|\bar{N}|}$$

unde  $P(u,v)$  este funcția de perturbare. Normala la suprafața perturbată este:

$$N' = Q'^u \times Q'^v = N + D$$

Funcția de perturbare poate fi orice funcție derivabilă,  $f(u,v)$ . Pentru texturile fără o definiție matematică, Blinn a sugerat utilizarea unui tablou bi-dimensional, indexat prin parametri  $u$  și  $v$ . Valorile intermediare se obțin prin interpolare bi-liniară iar derivatele parțiale prin diferențe finite.

Metoda se integrează bine într-un model de iluminare la nivel de fragment care folosește normala obținută prin interpolarea normalelor din vârfurile primitivei geometrice.

#### 2.3.3.2.2 *Hărți de normale (Normal maps)*

Ideea lui Blinn a fost extinsă, pe de o parte pentru a se ușura definirea detaliilor fizice ale suprafețelor, pe de alta parte pentru obținerea vitezei necesare sintezei imaginilor în timp real. Astfel, în locul funcției de perturbare a normalei se folosește o textură din care se iau normalele utilizate în calculul de iluminare la nivel de pixel. O astfel de textura se numește «hartă de normale».

Un vector normală este un vector 3D orientat de la suprafața înspre exterior. În cazul în care iluminarea este realizată în spațiul obiect, normalele stocate în textură corespunzătoare trebuie să fie reprezentate tot în spațiul obiect. Avantajul evident al acestei abordări este ca în momentul când sunt efectuate calculele de iluminare la nivel de pixel, normala poate fi citita din textura și utilizată direct, eventual după o renormalizare.

Această tehnică este folosită în mod uzual pentru simplificarea geometriei suprafețelor astfel încât să fie posibilă redarea lor în timp real. Astfel, un model poligonal extrem de complex poate fi înlocuit cu un model compus dintr-un număr rezonabil de poligoane și o hartă textură de normale care să îi adauge acestuia detaliile care s-au pierdut în urma simplificării.

Simplificarea obiectelor poate fi realizată în orice aplicație de modelare 3D, utilizând un instrument corespunzător. În general, aceasta reducere se face prin concatenarea poligoanelor vecine care au unghiul dintre planele lor sub un anumit prag, și prin concatenarea punctelor aflate la distanțe relative sub un alt prag specificat. Generarea texturii care memorează normalele se face pe baza diferențelor între obiectul redus și obiectul original, utilizând trasarea de raze pentru a deduce normalele peste obiectul redus.

#### **2.3.4 SUNET**

Pe măsură ce posibilitățile dispozitivelor de redare a sunetelor pe calculator și API-urile acestora s-au îmbunătățit, sunetul 3D a început să aibă un rol din ce în ce mai important în simularea lumilor virtuale.

Crearea unei experiențe cât mai imersive pentru utilizator în lumea virtuală nu înseamnă numai o reprezentare vizuală cât mai reală ci trebuie ținut cont și de elementele de sunet ale mediului virtual.

Implementate corect, sunetele ambientale, muzica și alte efecte audio pot duce simularea virtuală la un cu totul alt nivel.

#### **2.3.4.1 Efecte audio 3D**

Deși poziționarea în spațiul 3D este un factor important pentru modul în care utilizatorul percepe sunetele, aceasta singură nu este suficientă. Elementele mediului înconjurător pot modifica sunetul emis de o sursă înainte de a fi perceput de un personaj virtual existent în spațiul 3D. Principalele efecte acustice produse de interacțiunea sunetului cu mediul sunt:

- reverberația
- obstrucția
- ocluzia

Reverberația este fenomenul reflectării multiple a sunetului atunci când acesta întâlnește o suprafață.

Un factor de care trebuie ținut cont în simularea reverberației este materialul suprafeței de care se “lovește” unda sonoră, în funcție de acesta absorbția va fi mai mare sau mai mică rezultând o reverberație mai slabă sau mai puternică.

De asemenea, reverberația este influențată și de dimensiunea și forma incintei în care se propagă sunetul. Astfel, o incintă mai mare va face ca timpii de reverberație să fie mai lungi iar dacă avem suprafețe neregulate, cu multe muchii, acestea vor determina creșterea consistenței reverberației.

Reverberația ca fenomen acustic ne înconjoară în permanență în realitate și oferă foarte multe informații despre mediul înconjurător. De aceea, prezența acestui efect în simularea lumii virtuale contribuie semnificativ la realismul acesteia.

De asemenea, reverberația intervine și în cazul obstrucției acustice.

Obstrucția este fenomenul ce se produce atunci când între utilizator și sursa care emite sunetul se află poziționat direct un obstacol. Dacă se dorește ca simularea să fie cât mai apropiată de realitate, trebuie implementată și obstrucția sunetelor. În acest caz, undele sonore nu mai ajung direct la ascultător, dar totuși el le poate percepe într-o formă modificată, datorită reverberației.

Fenomenul de reverberație, care permite reflectarea sunetului, face posibil ca, prin reflexii multiple, acesta să ocolească obstacolul și, deși a pierdut din intensitate, să ajungă la utilizator.

Pe lângă aceste două fenomene, într-un spațiu real poate să apară și fenomenul de ocluzie acustică. Să luăm situația de mai sus, în care sursa de sunet este despărțită de utilizator printr-un obstacol, numai că de această dată obstacolul izolează perfect sunetul de utilizator, aceștia aflându-se în două încăperi diferite, neexistând nici o posibilitate ca sunetul să "ocolească" obstacolul (ca exemplu, putem considera două camere separate de un zid).

Astfel, dacă obstacolul este suficient de gros și izolant (de exemplu, piatra) sunetul nu va ajunge deloc la utilizator. Pe de altă parte, dacă obstacolul este alcătuit din lemn și este subțire sunetul va trece prin acesta atenuat, în funcție de grosimea obstacolului.

Fenomenul de atenuare a sunetului la trecerea printr-o suprafață de separație a două încăperi sau a două medii diferite se numește ocluzie.

Toate aceste fenomene sunt elemente importante în modelarea cât mai reală a lumilor virtuale. Pentru a obține aceste efecte, au fost dezvoltate diferite API-uri pentru efecte audio, cum ar fi Microsoft DirectSound și DirectSound3D, Aureal A3D și Creative Labs EAX.

Mai nou, Creative Labs a dezvoltat OpenAL, care este echivalentul pentru implementarea sunetului al OpenGL-ului.

#### **2.3.4.2 Formate Audio**

Există numeroase formate pentru fișierele ce conțin informațiile audio. În continuare sunt enumerate câteva dintre cele mai uzuale formate folosite în prezent.

##### **WAV**

Acesta este probabil formatul cel mai documentat folosit în prezent. Formatul .WAV este comparabil cu formatul .BMP pentru imagini în sensul că amândouă au o calitate foarte bună a informației pe care o conțin dar acest lucru este direct proporțional cu dimensiunea pe care o ocupă pe disc.

Deoarece aceste dimensiuni pot ajunge destul de mari este recomandat ca acest format să fie folosit pentru sunete și efecte audio de scurtă durată.

##### **MP3**

Este unul dintre cele mai populare formate audio. Acest format beneficiază de o tehnologie de compresie foarte bună cu pierderi minime de calitate. Un alt avantaj al acestui format ar fi ca este bine documentat, cu multe exemple online.

Acest format este recomandat a fi folosit pentru dialog și muzica ambientală.

## **OGG**

Acest format a început recent să capteze din ce în ce mai multă atenție. Deși încă nu este la fel de popular ca formatul MP3, recuperează rapid terenul pierdut. Față de MP3 are avantajul că se obține pentru sunete o compresie și o calitate mai bună.

### **2.3.5 ANIMAȚIA PERSONAJELOR**

În cadrul acestui sub-capitol sunt enumerate și prezentate pe scurt principalele metode de modelare și animare a personajelor (“characters”).

#### **2.3.5.1 Modelarea personajelor**

Metodele de modelare a personajelor au devenit din ce în ce mai complexe și reușesc să creeze modele foarte apropiate de personajele reale, dar nu întotdeauna modelul cel mai realist este și cel mai bun în orice aplicație. Alegerea uneia dintre metodele de modelare depinde de necesitățile aplicației.

##### *2.3.5.1.1 Plasa poligonală*

Acesta a fost primul mod de reprezentare a personajelor, direct inspirat din tehnicile de animație folosite în plan. Personajul este reprezentat printr-o plasă poligonală 3D. Animația personajului constă în deformarea plasei, fiecărui cadru imagine corespunzându-i o “poză” (ipostază) a personajului. Numărul de vârfuri ale plasei și conectivitatea lor nu se modifică de la o poză la alta, ci doar pozițiile vârfurilor.

##### *2.3.5.1.2 Reprezentari bazate pe schelet*

Scheletul unui personaj 3D animat este similar scheletului uman. Acesta constă din oase și articulații și poate fi folosit ca un mecanism de control pentru deformarea straturilor de deasupra scheletului: mușchi și piele. Scheletul este definit ca o ierarhie de oase sau articulații. El este reprezentat printr-un arbore ale cărui noduri corespund oaselor sau articulațiilor.

Unele aplicații care permit definirea de personaje reprezintă scheletul prin oase, altele prin articulații. Fiecare os/articulație are asociată o matrice de transformare (rotație și translație) care definește mișcarea articulației/osului. Transformarea totală a unui nod al ierarhiei schelet este produsul dintre transformarea nodului părinte și propria transformare. Animația scheletului constă în aplicarea de transformări, care se modifică în timp, asupra nodurilor.

Deoarece modelarea unui schelet real este destul de dificilă, în practică se folosește un schelet simplificat, cu câteva zeci de oase și articulații, cu unul, două sau trei grade de libertate, cărora uneori le sunt ignorate limitele.

Un avantaj al reprezentării prin schelet este faptul ca permite animatorului să controleze doar acele caracteristici ale modelului care se mișca independent. Animația este mult mai simplă de definit, rezumându-se la mișcarea oaselor, în loc de mișcarea fiecărui vârf ( în cazul reprezentării prin plasă poligonală).

### **2.3.5.2 Animația personajelor**

Există mai multe posibilități de a realiza animația personajelor iar în majoritatea cazurilor alegerea unei metode de animație este puternic dependentă de necesitățile aplicației țintă.

#### *2.3.5.2.1 Interpolarea între cadre cheie*

Metoda constă în realizarea de către animatori a unor schițe pentru personajul ce se dorește a fi animat. Schițele reprezintă imagini ale personajului la anumite momente importante. Ele se numesc “cadre cheie”. Apoi, se completează animația cu cadre intermediare între cadrele cheie succesive.

Sunt utilizate în acest scop diferite metode de interpolare. Pentru interpolarea pozițiilor se folosesc de regulă interpolarea liniară și cea pătratică. Pentru interpolarea orientării (între rotațiile corespunzătoare din două cadre cheie) se folosește interpolarea sferică, folosind reprezentarea rotațiilor prin quaternioni (Slerp - interpolarea sferică liniară și Squad - interpolarea sferică cubică).

Interpolarea între cadre cheie a fost folosită la început pentru animația reprezentării prin plasă poligonală, caz în care se aplică vârfurilor plasei. Ulterior ea a fost asociată și altor tehnici de animație, de exemplu, celor bazate pe animația scheletului, caz în care cadrele cheie reprezintă ipostaze ale scheletului iar interpolarea se folosește pentru generarea matricilor de transformare ale oaselor.

#### 2.3.5.2.2 Interpolarea formei

Interpolarea formei (Shape Interpolation) este o tehnică folosită în special în animația facială. Se definește un set de forme “cheie”. Coordonatele punctelor de pe suprafața formei sunt calculate ca o combinație liniară a coordonatelor punctelor din formele cheie, pe baza unor ponderi asociate punctelor:

$$S = \sum_k w_k \cdot S_k$$

O variație a acestei tehnici folosește o formă ca bază, celelalte fiind specificate doar ca diferențe față de această bază:

$$S = S_0 + \sum_k w_k \cdot (S_k - S_0)$$

Ponderile cu care se combină aceste date variază în timp.

Metoda are o aplicabilitate limitată, deoarece se utilizează pentru combinare forme rigide.

#### 2.3.5.2.3 Forward Kinematics

Este o metodă folosită pentru specificarea animației unei structuri articulate, cum ar fi scheletul unui personaj. Metoda presupune specificarea explicită de către animator a mișcărilor fiecărei părți a structurii articulate. Animatorul trebuie să definească mișcările obiectelor rigide din structură (oasele) începând cu nodul rădăcină al ierarhiei, coborând în ierarhie până la fiecare terminație a structurii.

#### 2.3.5.2.4 Inverse Kinematics

În această metodă animatorul trebuie să specifice numai poziția (pozițiile) în care trebuie să ajungă terminația (sau terminațiile) structurii. Prin cinematica inversă se deduc mișcările tuturor părților structurii articulate astfel încât să se obțină “poza” dorită.

Problema este că pe măsura ce structurile devin din ce în ce mai complexe, soluția cinematicii inverse devine din ce în ce mai greu de găsit. Există o rezolvare analitică a problemei doar pentru lanțuri foarte scurte. În cazul unei ierarhii precum cea a scheletului personajelor umanoide, soluția nu poate fi găsită decât folosind tehnici iterative. Din acest motiv metoda este lentă.



#### *2.3.5.2.5 Inverse Kinetics*

Înrudită cu cinematica inversă, permite ierarhizarea scopurilor pe care trebuie să le atingă personajul de animat și specificarea unui criteriu în funcție de care trebuie optimizată soluția. Criteriul este exprimat ca o funcție de starea personajului. În plus față de cinematica inversă, această tehnică ține cont de distribuția masei, obținând astfel rezultate mai bune, mai repede. Deși are aceeași complexitate, folosind metode asemănătoare, cinetica inversă converge mai repede în cazul încercării de a poziționa centrul de masă.

Influența fiecărei articulații asupra posturii căutate este estimată în funcție de masa pe care o suportă în starea curentă a sistemului. Rotația oaselor în articulații, este pusă în legătură cu translația unui end-effector la fel ca în cazul cinematicii inverse.

#### *2.3.5.2.6 Forward Dynamics*

Metodele de cinematica ignoră o mare parte din influențele mediului. Nu țin cont de masa obiectului de animat sau de forțele care-i sunt aplicate.

Această metodă calculează pozițiile folosind distribuția masei personajului și forțele care se aplică. În cazul jocurilor se pot folosi legile mecanicii Newtoniene. Știind forțele, putem calcula accelerația, care apoi este integrată și se calculează poziția. Deși robotica a dezvoltat algoritmi eficienți pentru structuri cu multe grade de libertate, metoda este greu de folosit. Forțele nu sunt niște parametri sugestivi pentru un animator. Controlul unui personaj astfel simulat este foarte dificil.

#### *2.3.5.2.7 Inverse Dynamics*

În aceasta metodă, pornind de la poziția ce se dorește a fi obținută se încearcă aflarea forțelor care trebuie aplicate.

Rezultatele obținute sunt în general aproximative. Dinamica inversă face anumite presupuneri: că nu există forțe de frecare la articulații, că oasele sunt rigide și masa lor e concentrată în centrul de masă și nu există frecare cu aerul. Are aplicații în biomecanică. Încearcă să afle forțele produse de mușchi în timpul mișcării și în acest scop presupune ca mușchii nu-și anulează reciproc acțiunile.

#### *2.3.5.2.8 Motion Capture (captura mișcării)*

“Captura mișcării” are drept scop înregistrarea mișcărilor unui personaj real, cu ajutorul unor echipamente, pentru a fi apoi folosită în animația personajelor

sintetizate. Informația capturată variază de la poziția diverselor părți ale corpului până la deformări ale feței sau mușchilor.

Există mai multe tipuri de echipamente de captură a mișcării: magnetice, optice, electro-mecanice, opto-mecanice ce folosesc pentru înregistrarea mișcării senzori și markeri ce sunt plasați pe corp.

Animațiile produse prin captura mișcării nu sunt decât varianta modernă a animației bazată pe cadre cheie (cadrele cheie nu sunt desenate ci înregistrate) și ele pot servi ca date de foarte bună calitate și care pot fi obținute rapid, pentru alte metode.

#### *2.3.5.2.9 Motion Warping (deformarea mișcării)*

Mișcarea poate fi văzută ca un set de curbe care descriu valorile parametrilor modelului animat în funcție de timp. Motion Warping pleacă de la convingerea că se pot modifica anumite traiectorii, pentru a respecta unele constrângeri, fără a afecta credibilitatea animației. Principala aplicație a fost modificarea datelor provenite din captura mișcării pentru a putea fi folosite în situații specifice. Metoda păstrează într-adevar aspectul natural al mișcării, dar curbele sunt dificil de înțeles și editat.

#### *2.3.5.2.10 Motion Blending (combinarea mișcărilor)*

«Motion Blending» a pornit de la ideea ca se pot combina mai multe animații pentru a obține o animație diferită. De exemplu, se poate combina o animație care reprezintă un om mergând furios, cu o animație a aceluiași om care aleargă, pentru a obține un om furios care aleargă.

Se pot aplica anumite constrângeri fizice. De exemplu, înălțimea la care ajunge mâna în cazul unei lovituri.

Metodele de combinare a mișcărilor par a fi promițătoare pentru sistemele de realitate virtuală. Deși nu respectă în totalitate legile fizicii, ele produc în general mișcări care arată natural, iar unele dintre ele pot fi controlate cu parametri de nivel înalt. Astfel de parametrii pot fi folosiți pentru modelarea relației agenților cu mediul virtual.

### **2.3.6 SISTEME DE PARTICULE**

Termenul de **sistem de particule** se referă la o tehnică folosită în grafica pe calculator, prin care se pot simula anumite fenomene sau obiecte neregulate cum sunt: focul, exploziile, fumul, apa, scânteile, frunzele, nori, ceață, zăpada, praful sau efecte vizuale abstracte cum ar fi o urmă strălucitoare, etc. Astfel de obiecte sau fenomene sunt greu de reprezentat și redat în imagini prin tehnici convenționale.

De regulă, poziționarea sistemului de particule în spațiul 3D și dinamica acestuia sunt controlate de către un **emițător**.

Emițătorului îi este atașat un set de parametri ce caracterizează comportamentul particulelor: rata de generare a acestora (câte particule sunt generate într-o unitate de timp) , vectorul inițial al vitezei particulelor (în ce direcție se vor mișca particulele când acestea vor fi emise) , durata de viață a particulelor, variația culorii particulelor în timpul vieții lor , etc.

Atributele unei particule individuale sunt : poziția, viteza, mărimea, culoarea, transparența, forma, durata de viață.

În mod normal, toți acești parametri nu au valori regulate, acestea putând varia de la o particulă la alta, fiind aleatorii între anumite limite.

O comportare caracteristică pentru un sistem de particule poate fi împărțită în două etape:

- etapa de actualizare a stării sistemului
- etapa de redare a sistemului (crearea imaginii curente a sistemului)

**În timpul etapei de actualizare** sunt generate noi particule și se actualizează starea celor existente în sistem. Fiecare nouă particulă este creată într-o poziție 3D specifică, bazată pe poziția emițătorului și pe zona de generare specificată, iar atributele ei (viteza, culoarea) sunt inițializate în funcție de parametrii emițătorului.

Particulele existente în sistem a căror durata de viață a expirat sunt eliminate din simulare.

Starea celorlalte este modificată pe baza unui mecanism de simulare, care poate să fie foarte simplu, constând în calculul noii poziții folosind viteza de deplasare, sau unul complicat în care se fac calcule ale traiectoriei particulei corecte din punct de vedere fizic, care iau în considerare acțiunea unor forțe externe (cum ar fi de exemplu gravitația).

În mod normal, tot în această etapă ar trebui verificată detecția coliziunilor particulelor cu obiectele 3D din sistem (sau din spațiul virtual), astfel încât particulele să ricoșeze la întâlnirea acestora. Totuși verificarea detecției coliziunilor în cazul sistemelor de particule este rar utilizată deoarece este foarte costisitoare din punct de vedere computațional.

Un sistem de particule are propriile reguli, care sunt aplicate fiecărei particule în parte. Uzual aceste reguli conțin și o metoda de interpolare a valorilor atributelor particulei de-a lungul vieții sale. De exemplu, multe sisteme fac ca particulele să dispară ușor prin interpolarea valorii alfa (opacitatea) de-a lungul vieții particulei (“fading out”).

Dupa ce etapa de actualizare s-a incheiat, fiecare particulă este redată în imagine, de obicei sub forma unui billboard. Aceasta nu este însă singura modalitate de a reprezenta o particulă : se pot folosi reprezentări de la cea mai simplă, un singur punct, până la reprezentări de suprafețe (mesh-uri).

## **2.3.7 REDĂRI SPECIALE**

### **2.3.7.1 Redarea vegetației**

Absența redării cu mare acuratețe a vegetației în aplicațiile de timp real, cum ar fi și simulările virtuale, poate fi atribuită în mod direct cantității imense de informație geometrică necesară pentru a modela corespunzător vegetația.

Iarba și alte forme de vegetație naturală sunt frecvent întâlnite în simulările spațiilor exterioare. Deși, de exemplu, forma unui singur fir de iarba nu este complexă și aceasta poate fi reprezentată cu ajutorul câtorva poligoane, firele de iarba nu apar într-un număr mic, acestea întinzându-se de obicei pe suprafețe foarte mari. Astfel, reprezentarea și iluminarea vegetației devin greu de realizat din punct de vedere practic.

Mai mult, vegetația poate avea o variație structurală foarte mare ceea ce duce la apariția de proprietăți caracteristice de reflexie a luminii.

Redarea vegetației prin modelarea sa geometrică este posibilă dar nu este o metodă viabilă în practică datorită cantității de informație foarte mare care necesită putere computațională ridicată.

O soluție constă în simularea vegetației folosind texturi.

Utilizarea de texturi statice produce rezultate destul de apropiate de realitate, dar exista două probleme importante în utilizarea acestei soluții, care scad din realismul simulării :

- texturile folosite fiind statice, nu se modifică în funcție de poziția camerei și de iluminare; aceasta înseamnă că, de exemplu, un petic de iarba va arăta la fel din orice poziție este privit, umbrele vor rămâne fixe, spoturile speculare nu se vor modifica, etc; astfel, natura statică a texturii scade foarte mult realismul reprezentării grafice a vegetației;
- texturile fiind imagini bidimensionale, sunt mapate pe un poligon plat astfel încât realismul 3D al mediului scade simțitor; nu se pot obține siluete și efectul interacțiunii vegetației cu alte obiecte din mediul înconjurător dispore.

O altă metodă, care reprezintă o evoluție a folosirii texturilor statice, este folosirea de billboard-uri. Transparența texturii aplicate pe suprafața billboard-ului face posibilă reprezentarea de obiecte complexe, cum ar fi plantele și copacii. Metoda însă nu permite realizarea unei iluminări realiste.

Alte abordări se bazează pe utilizarea de texturi volumetrice ([DEN04]).

Totodată, evoluția hardware-ului grafic a creat o nouă provocare în privința redării vegetației folosind modele geometrice și iluminare realistă([CDN04]).

În practica se folosește o schemă hibridă, bazată pe un model cu trei nivele:

- pentru vegetația care se află aproape de poziția camerei se folosește o reprezentare geometrică detaliată;
- pentru vegetația aflată la distanță medie se folosesc texturi volumetrice;
- pentru vegetația aflată la distanță mare, pentru care oricum detaliile nu s-ar putea distinge, se folosește o schemă de texturare statică.

Sinteza imaginilor realiste a terenurilor acoperite cu vegetație este și în prezent o provocare și o problema importantă în grafica pe calculator.

### **2.3.7.2 Redarea apei**

Simularea apei constituie un subiect asupra căruia se fac cercetări de un timp îndelungat. Simularea realistă a fluidelor s-a dovedit a fi foarte dificilă datorită complexității inerente a proceselor fizice ce intervin în comportarea acestora.

Cele mai realiste soluții folosesc un sistem de tip ray-tracing și de ecuații diferențiale complexe pentru a aproxima comportamentul optic și dinamica

mișcării apei. Din nefericire aceste abordări de o acuratețe fizică mare pot necesita multe ore de procesare numai pentru un singur cadru.

O dată cu dezvoltarea continuă a hardware-ului programabil pentru grafica 3D, mulți algoritmi care înainte erau folosiți doar pentru procesare off-line sunt acum viabili pentru folosirea în timp real.

Ca de obicei, atunci când se implementează efecte complexe în contexte interactive, trebuie ajuns la un compromis între calitatea reprezentării și performanța computațională necesară pentru a face această reprezentare. Astfel, foarte multe proprietăți vizuale ale apei pot fi aproximate prin tehnici eficiente, care deși nu sunt corecte din punct de vedere fizic par convingătoare pentru ochiul uman.

La fel ca majoritatea efectelor fizice, simularea apei este o combinație a doi pași distincți :

- **Generarea fizică** a apei, denumită și modelul undei, simulează mișcarea apei sub influența unor forțe externe și interne; acest model este cel care face ca apa să se miște, să își schimbe forma și să interacționeze cu mediul înconjurător. Din punct de vedere matematic, sunt posibile mai multe abordări. Cele mai simple aplică o perturbare procedurală, cel mai adesea sub forma unei perturbări Perlin, simulând o mișcare ce pare aleatoare. Metodele care încearcă să reproducă mișcarea reală a apei folosesc ecuații diferențiale, aproximarea legilor fizice ale hidrostatiei și hidrodinamicii și permit interacțiunea dintre utilizator și apă.
- **Vizualizarea** apei pe baza datelor calculate la pasul anterior. Pentru un realism cât mai ridicat trebuie să se țină iarăși cont de legile fizicii : o suprafață de apă are interacțiuni complexe cu lumina din mediul înconjurător și modelarea acestor proprietăți optice este foarte importantă pentru obținerea unui efect vizual cât mai real.

Ambii pași pot beneficia de puterea de calcul a GPU, dar mai ales al doilea. Simularea optica a suprafeței apei necesită operații complexe de calcul la nivel de pixel. Hardware-ul modern cu posibilități de programabilitate atât la nivel de vârf cât și la nivel de pixel este potrivit pentru această sarcină.

## **2.4 TEHNOLOGII SERVER-SIDE**

În proiectarea protocolului de comunicație dintre server și utilizatori, într-o aplicație spațiu virtual care se execută în cadrul unei rețele pe mai multe calculatoare, se poate opta pentru abordarea “peer-to-peer” sau abordarea “client/server”.

### **2.4.1 ABORDAREA PEER-TO-PEER**

În abordarea peer-to-peer, două sau mai multe calculatoare comunică între ele, fiecare rulând o copie a aplicației și schimbând informații despre starea aplicației.

Aceasta abordare, deși conduce la un trafic redus pe rețea, implică alte probleme cum ar fi, de exemplu, tratarea traficului pierdut. Aceasta este mult mai importantă în momentul în care rulează mai mult de o copie a aplicației, deoarece nu se mai știe cine are datele corecte despre starea aplicației și cine nu. Diferențele între stările aplicației la diferiți utilizatori pot crea probleme, deoarece toate copiile aplicației trebuie să se sincronizeze între ele.

În plus, fiecare instanță a aplicației trebuie să aștepte intrările de la celelalte instanțe înainte de a putea simula următorul cadru. Doom 3 [<http://www.doom3.com/>] a fost unul dintre jocurile care a folosit această abordare și existau momente în care jocul “îngheța” în așteptarea stării de la mașinile celorlalți jucători.

### **2.4.2 ABORDAREA CLIENT/SERVER**

În abordarea client/server, în sistem există doi actori :

- mașina server
- clienții

Mașina server este cea care rulează aplicația și transmite către toți clienții starea curentă a aplicației, pe baza căreia fiecare client își actualizează imaginea spațiului virtual. Clienții afișează scena spațiului virtual așa cum le-a fost transmisă de către server și redau sunetele pe care server-ul “le spune” să le redea.

Efectiv, clienții se comportă ca niște terminale care transmit intrările către server și îl lasă pe acesta să se ocupe de aproape tot. În practică, acest mod de împărțire a

responsabilităților nu este un lucru chiar așa de rau cum pare la prima vedere, ocazional server-ul putând sa ceară unui client să îndeplinească și alte funcționalități ce ar putea sa cadă în seama sa.

Pot exista foarte multe variații în abordarea client/server însă ideile de bază sunt cele prezentate în continuare.

#### **2.4.2.1 Probleme ce trebuie tratate în realizarea proiectării sistemului client/server**

În proiectarea unui protocol multi-utilizator, pot apărea următoarele probleme, de care trebuie să se țină cont :

- tehnologia folosită : TCP sau UDP
- încărcarea pachetelor
- frecvența pachetelor
- ordinea pachetelor și pachetele lipsă (în cazul UDP)
- clienți care încearcă să păcălească sistemul

#### **Tehnologia folosită: TCP sau UDP**

UDP (User Datagram protocol) și TCP (Transmission Control Protocol) sunt două protocoale ce funcționează la nivelul Transport al modelului OSI.

Diferența dintre UDP și TCP este aceea ca TCP garantează livrarea pachetelor la destinație în ordinea în care au fost transmise pe când UDP nu garantează livrarea lor și cu atât mai puțin ordinea în care acestea sosesc.

Deși la o prima vedere TCP pare sa fie superior fata de UDP, avantajele oferite de protocolul TCP au un cost: overhead-ul (timpul suplimentar consumat cu stabilirea conexiunii) și latența.

Astfel, în protocolul TCP, expeditorul (ca sa fie sigur ca pachetele ajung la destinație) așteaptă să primească un ACK (Acknowledgement) din partea destinatarului pentru fiecare pachet pe care îl trimite. Dacă nu primește nici un ACK atunci încetează să mai trimită pachete noi și le retrimite pe cele presupuse pierdute, continuând sa facă acest lucru până când destinația răspunde.

Problema este întârzierea dintre momentul în care expeditorul încetează să mai trimită pachete noi și momentul în care pachetul ajunge într-adevăr la destinație.

Această întârziere poate ajunge uneori la nivel de secunde, ceea ce nu ar fi un motiv de îngrijorare dacă s-ar descărca o pagina sau un fișier, însă în cazul unui



pachet dintr-un joc apar probleme, deoarece aceasta împiedică actualizarea stării jocului la destinație.

#### **2.4.2.2 Încărcarea pachetelor**

Trebuie acordată o atenție foarte mare la selectarea datelor transmise, astfel încât să fie transmise numai acele date care sunt strict necesare. Cu cât este mai mare pachetul trimis, cu atât se încarcă mai mult transmisia ajungându-se astfel la latențe mari dacă conexiunea folosită nu are o lățime de bandă suficient de mare.

Trebuie ținut cont și de faptul că un pachet ce se poate trimite are o dimensiune maximă numită MTU (maximal transmission unit). De obicei, în cazul protocoalelor multiplayer se definește o dimensiune maximă a pachetului iar atunci când se depășește această dimensiune, datele se împart în mai multe pachete, transmițându-se mai întâi datele care trebuie să ajungă cât mai repede.

O altă tehnică prin care se poate reduce încărcarea pachetelor este eliminarea mesajelor text. Dacă serverul trimite foarte multe mesaje/comenzi text prestabilite este logic ca acestea să fie pre-încărcate la client și acesta să primească de la server doar identificatorul mesajului în loc de mesaj. Această tehnică poate reduce traficul de mesaje în mod considerabil.

Trimiterea doar a obiectelor/datelor care au relevanță pentru scena vizualizată de utilizator este de asemenea indicată. Clientul nu are nevoie să cunoască informații despre acțiuni și evenimente care nu se află în raza sa vizuală sau auditivă. Oricum acestea nu vor fi redată sau auzite de client. Acest lucru ar putea să fie în cadrul unei scene exterioare sau al unui simulator în spațiu să nu fie tocmai indicat, însă și aici se pot lua decizii de implementare care să compenseze, cum ar fi, de exemplu, implementarea nivelului de detaliu.

Nu este necesar ca la client să fie transmise rezultatele execuției efectelor speciale, acest lucru consumând lățimea de bandă atât a serverului cât și a clientului. Este de ajuns, de exemplu, ca serverul să spună “aici se întâmplă o explozie” și clientul face restul afișând efectul corespunzător pe display.

O altă tehnică folosită este *compresia bazată pe modificare*. Ideea este să se transmită numai ceea ce s-a modificat între două cadre succesive. Se poate păstra o copie a ceea ce s-a transmis ultima oară, bineînțeles la nivel de obiect, și pentru cadrul curent să se transmită numai modificările.

Ca ultimă idee, se pot implementa tehnici de compresie a datelor din interiorul pachetelor trimise pentru a micșora dimensiunea pachetelor, însă de cele mai multe ori timpul consumat atât pe server cât și pe client pentru compresie/decompresie este mai mare decât trimiterea unui pachet necompresat.

#### **2.4.2.3 Frecvența pachetelor**

Deși motorul grafic este în stare să redea scena la un număr mare de cadre pe secundă, dacă infrastructura rețelei nu poate să susțină o rată de trimitere/sosire a pachetelor suficient de mare, acest lucru nu este un avantaj.

De exemplu, serverul de *Quake* funcționează la 10 cadre pe secundă, transmițând datele la această rată. Funcționează la 10 cadre pe secundă din două motive:

- din cauza dimensiunii datelor pe care le procesează pentru fiecare client
- deoarece este o rată a pachetelor care poate fi susținută ușor

#### **2.4.2.4 Ordinea pachetelor și pachetele lipsă**

Modul de abordare este de a trata un simptom și de a-l ignora pe celălalt. Dacă pachetele sunt numerotate, clientul va ști când a primit un pachet ce nu a sosit în ordinea corectă. Cea mai simplă soluție este să îl ignore și să îl trateze ca pe un pachet pierdut.

Dacă la client se păstrează o copie a ultimului pachet primit de la server, se poate compara pachetul ce sosește cu acesta pentru a vedea dacă un obiect a fost eliminat din scenă. Dacă acest lucru s-a întâmplat, se poate elimina și de către client obiectul imediat sau se poate stoca într-un buffer pentru a mai verifica cu încă câteva pachete ce vor sosi. Avantajul acestei tehnici este că nu va fi necesar niciodată să se trimită de către server clientului comanda de “remove” deoarece lipsa unui obiect din pachet garantează eliminarea acestuia și de la client. Chiar dacă au existat pachete lipsă, acest lucru nu contează deoarece până la urma clientul va primi un pachet și va vedea că obiectul nu mai există !

Tot aici se poate discuta despre mecanismul de predicție a acțiunilor clientului. În cazul în care apar pachete lipsă, atât din partea clientului cât și a serverului, în timpul până la primirea unui pachet corect, clientul trebuie să execute o acțiune pentru a nu da impresia de “înghețare”. Se poate prezice până la un anumit punct ce se va întâmpla în următoarele cadre la client (acesta poate avansa, sau executa o săritură, etc) și să se afișeze pe ecran efectele. Desigur acest mecanism

funcționează la nivel de fracțiuni de secundă și în cel mai bun caz la nivel de câteva secunde, dar de obicei este de ajuns pentru ca sistemul de trimitere/primire a pachetelor să își revină și clientul să primească din nou pachete de la server, moment în care clientul poate să corecteze acțiunile care au fost prevăzute greșit.

În cel mai bun caz, acțiunile au fost prevăzute corect și clientul nici nu va ști că s-au pierdut date, în cazul mediu se fac mici corecții imperceptibile, iar în cel mai rău caz starea clientului s-a modificat foarte mult din cauza lipsei pachetelor de la server.

#### **2.4.2.5 *Clienți care încearcă să păcălească sistemul***

Există câteva modalități de a trata această problemă, dar nu există o soluție care să funcționeze în toate cazurile, găsindu-se mereu noi modalități de a ocoli orice protecție.

Nu este plăcut pentru nimeni să joace un joc în care cineva este invincibil datorită faptului că trișează. Acest lucru se poate produce în multe feluri, de exemplu modificând clientul să nu arate zidurile, adăugând lumini sau elemente ce pot scoate în evidență poziția celorlalți jucători foarte ușor, modificând ținta astfel încât de fiecare dată să nimerească adversarul cu precizie absolută.

Toate aceste lucruri sunt modificări ce se întâmplă în partea de client a jocului și când sunt făcute așa cum trebuie sunt aproape de neobservat de către server.

S-ar putea, de exemplu, verifica precizia fiecărui jucător și eliminarea celor care au o precizie care depășește un anumit prag. În acest fel se pot pierde jucători foarte buni dar este puțin probabil ca cineva poate avea o precizie de peste 80% ([KES05]) tot timpul.

Integritatea clientului este aspectul important aici. Soluția adoptată de jocurile Quake1 și Quake3, aceea a unei mașini virtuale unde în locul încărcării jocului de către client acesta “construiește” și “compilează” jocul în urma instrucțiunilor primite de la server este un start bun.

Tot ceea ce poate face dezvoltatorul de jocuri / spații virtuale este să facă sarcina celui care vrea să păcălească sistemul cât mai dificilă.

## **2.5 CONCLUZII**

Principala caracteristică a aplicațiilor MMO este dată de faptul că oferă posibilitatea unui număr ridicat de utilizatori să se întâlnească și să interacționeze în cadrul aceleiași lumi virtuale. În funcție de scara aplicației MMO, acest număr poate să varieze de la câteva sute ajungând până chiar la zeci de mii de utilizatori ce accesează simultan spațiul virtual. Pentru a crea un grad de realism ridicat utilizatorilor ce folosesc spațiul virtual, aplicațiile MMO trebuie să ofere pe partea de client o redare 3D cât mai bogată din punct de vedere vizual și pentru a realiza acest lucru sunt utilizate motoare grafice 3D în timp real moderne.

În acest capitol au fost prezentate categoriile de aplicații MMO cele mai des întâlnite pe piață și au fost detaliate conceptele specifice folosite de acestea. De asemenea au fost trecute în revistă tehnologiile 3D majore care sunt folosite de motoarele grafice 3D în timp real.

Noțiunile descrise în acest capitol au fost prezentate și într-o carte la care autorul lucrării de față a fost co-autor ([MMA09a]).

## **3 ANALIZA SOLUȚIILOR ARHITECTURALE ACTUALE PENTRU SERVERE DE SPAȚII VIRTUALE 3D MMO**

### **3.1 FUNCȚIONALITĂȚI DE BAZĂ**

O arhitectură a unei aplicații MMO trebuie să ofere următoarele funcționalități de bază:

- ❖ Autentificare utilizatori
- ❖ Gestionarea drepturilor de acces ale utilizatorilor în lumea virtuală
- ❖ Consistența lumii virtuale
- ❖ Gestiunea ordinii evenimentelor
- ❖ Propagarea evenimentelor către entitățile spațiului virtual
- ❖ Stocarea securizată a caracterelor și posesiunilor acestora
- ❖ Planificarea operațiilor computaționale
- ❖ Efectuarea calculelor de actualizare a stării lumii virtuale
- ❖ Timpi de raspuns scăzuți
- ❖ Securitatea lumii virtuale pentru a preveni eventualele încercări de a trișa

Se pot astfel identifica în cadrul unei arhitecturi generice următoarele componente fundamentale:

- 1) Componenta de **Autentificare** : este responsabilă cu controlul accesului la spațiul virtual
- 2) Componenta de **Comunicare** : este responsabilă cu schimbul de mesaje și gestiunea evenimentelor. Principalele sarcini ale acestei componente sunt de a menține ordinea corectă a evenimentelor la nivelul întregului spațiu

virtual, de a menține o latență cât mai scăzută și de a asigura mecanisme de securitate pentru a preveni încercările de trișare.

- 3) Componenta de **Stocare** : asigură stocarea persistentă a datelor lumii virtuale pe termen lung. Datele ce se pot stoca pot fi împărțite la rândul lor în două categorii : permanente și temporare
- 4) Componenta **Computațională** : este responsabilă cu planificarea și execuția operațiilor ce țin de logica spațiului virtual. De asemenea, tot în cadrul acestei componente este necesară implementarea de mecanisme de securitate pentru a asigura corectitudinea operațiilor executate de utilizatori.

## **3.2 ZONARE**

Cum s-a menționat și anterior, un aspect important de luat în considerație atunci când se proiectează arhitectura este faptul ca aplicațiile MMO trebuie să simuleze lumi virtuale de întindere foarte mare și acest lucru face necesară existența unui număr ridicat de servere, uneori de ordinul sutelor, pentru a găzdui în mod complet mediul virtual. Astfel, a apărut necesitatea dezvoltării unor metode de a împărți spațiul virtual în mai multe subspații distincte (*[GLI08]*).

În prezent, exista două abordări distincte atunci când vorbim de împărțirea spațiului virtual în subspații.

### **3.2.1 ZONARE NETRANSPARENTĂ UTILIZATORULUI**

Cea mai simplă din punct de vedere al implementării este cea a împărțirii stricte pe zone care practic partiționează întreaga lume virtuală în diferite părți. Aceste părți sunt suficient de mici pentru a putea fi gestionate de un singur server.

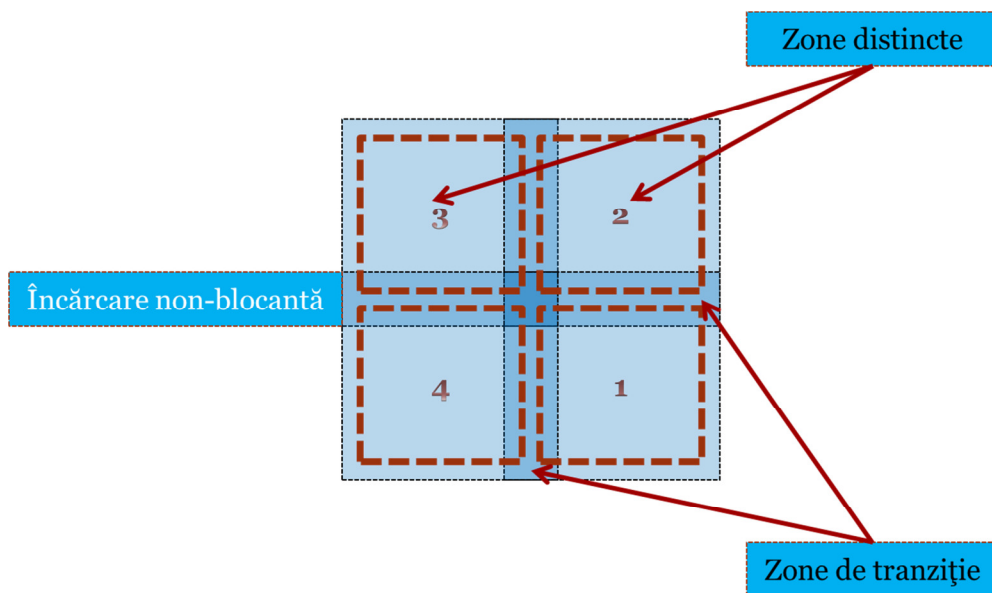
Granițele acestor zone sunt conectate între ele din punct de vedere logic și un utilizator se va conecta la serverul corespunzător atunci când trece dintr-o zonă în alta. Dezavantajul evident al acestei abordări este dat de faptul că există un timp de încărcare atunci când se schimbă zonele, acesta ducând la scăderea realismului simulării (de obicei, atunci când este vorba de un MMO dinamic, utilizatorii pot traversa un număr mare de zone într-un interval de timp relativ scurt pentru a îndeplini diverse obiective în cadrul spațiului virtual).

### 3.2.2 ZONARE TRANSPARENTĂ UTILIZATORULUI

Cea de-a doua abordare își propune crearea unei simulări “seamless” prin care utilizatorul să nu își dea seama de existența zonelor și a granițelor acestora. Acest lucru va crea utilizatorului senzația că se află într-o lume ‘continuă’ care funcționează fără ‘întreruperi’.

Deși nu este vizibilă explicit utilizatorului, separarea logica în zone a lumii virtuale există. La fel ca la abordarea precedentă, zonele importante care se învecinează sunt gestionate de servere diferite. În exemplul din figura de mai jos, atunci când utilizatorul iese din zona gestionată de serverul zonei 1, va exista o zonă de tranziție care este sincronizată cu zonele adiacente. Preluarea utilizatorului de către serverul zonei 2 nu se va efectua imediat ci doar atunci când acesta va traversa complet zona de tranziție.

Dezavantajul acestei abordări este complexitatea ridicată din punct de vedere al programării precum și posibilitatea apariției unui trafic crescut atunci când există multe evenimente în zonele de tranziție.



*Figura 3-1: Zonare fără „întreruperi” a lumii virtuale*

### **3.3 FOLOSIREA DE INSTANȚE**

Într-o aplicație MMO, o instanță reprezintă o zonă specială care își creează câte o copie a ei (instanțele sunt create și gestionate de obicei în cadrul arhitecturii spațiului virtual de către grupuri de servere dedicate acestora), de unde și denumirea de „instanță”, pentru un anumit grup de utilizatori care doresc să o acceseze. Astfel, acel grup de utilizatori accesează practic acea copie obținându-se astfel o izolare față de ceilalți utilizatori ai spațiului virtual deoarece conținutul instanței este vizibil numai în contextul grupului de utilizatori care a accesat-o.

De obicei, exista servere dedicate pentru instanțe, acestea reprezentând zone mai deosebite ale spațiului virtual, al căror design nu ar permite accesul unui grup numeros de utilizatori simultan.

Câteva dintre avantajele folosirii de instanțe ar fi :

- Reducerea operațiilor ce trebuie executate de serverele care simulează spațiul virtual global, acestea fiind transferate serverelor pentru zonele instanțelor.
- Reducerea traficului de rețea al utilizatorilor, acesta fiind practic izolat în cadrul instanței în care se află, obținându-se astfel un timp de răspuns scăzut pentru aceștia.
- Posibilitatea proiectării unor experiențe focalizate, având ca țintă numai un anumit număr de utilizatori
- Eliminarea competiției, care poate crea uneori mari probleme, asupra unor anumite resurse / obiecte datorită limitării numărului de utilizatori care pot accesa o instanță.

Cu toate aceste avantaje însă este de observat ca există și un dezavantaj al folosirii de instanțe și anume cel al scăderii realismului simulării spațiului virtual datorită izolării utilizatorilor din instanță față de restul spațiului virtual.



## **3.4 ARHITECTURI TRADIȚIONALE CLIENT-SERVER**

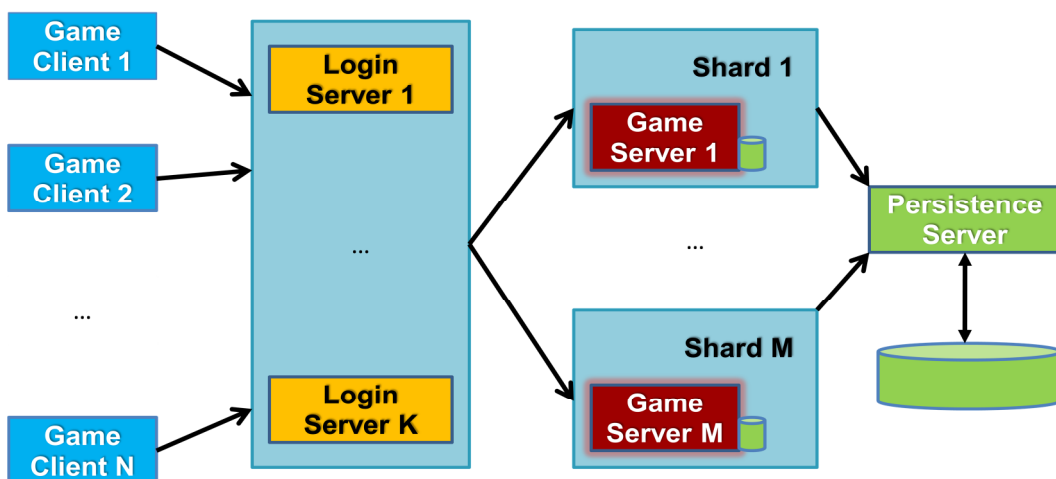
### **3.4.1 DESCRIERE GENERALĂ**

Marea majoritate a aplicațiilor de tip MMO folosesc o arhitectură de tipul Client-Server. Clienții accesează spațiul virtual prin intermediul unui *server de conectare* care apoi îi direcționează către un *shard server*.

*Shard-urile* reprezintă versiuni independente ale aceleiași lumi virtuale și se folosesc pentru a îmbunătăți scalabilitatea aplicației. *Shard-urile* nu sunt sincronizate și utilizatorii de pe un *shard* nu pot interacționa cu utilizatorii de pe alt *shard*.

Deși această soluție permite un anumit grad de scalabilitate al aplicațiilor MMO, ea limitează gradul de realism împiedicând interacțiunea unui număr mare de utilizatori.

Unele aplicații MMO, cum ar fi, de exemplu, *Eve Online* ([EVE11]) încearcă găsirea de noi soluții, care să nu folosească *shard-uri* și să permită interacțiunea **tuturor** utilizatorilor lumii virtuale. Aceste încercări sunt încă în faza incipientă și momentan depind foarte mult de particularitățile lumii virtuale.



*Figura 3-2: Arhitectura clasică MMO client-server*

Arhitectura client-server este în general preferată deoarece un sistem centralizat oferă în mod natural următoarele avantaje:

- consistența lumii virtuale
- securitate ridicată
- evitarea problemelor de sincronizare
- implementare mai simplă
- model de afaceri viabil

### **3.4.2 INTERACȚIUNEA DINTRE APLICAȚIILE CLIENT ȘI SERVER**

Serverul și clientul au ambele o arhitectură bazată pe nivele, nivelul cel mai de jos ocupându-se de comunicarea prin Internet folosind protocoale de comunicare standard. Din motive de securitate, următorul nivel este unul care ascunde formatul intern al datelor transmise folosind mecanisme de criptare/decriptare.

Ultimul nivel este cel aplicație și el reprezintă nivelul arhitectural propriu-zis al spațiului virtual 3D MMO, fiind cel mai complex nivel al aplicației.

Pentru a face transparentă nivelului aplicație complexitatea nivelului de criptare, se folosește un nivel intermediar de mesaje între acestea două.

Din punct de vedere conceptual, aplicația client interacționează cu serverul într-o manieră specifică șablonului de proiectare *publisher / subscriber*. Datorită naturii aplicațiilor MMO, serverul trebuie să comunice cu un număr ridicat de clienți, astfel pentru a minimiza traficul, serverul nu îi anunța pe aceștia imediat ce are loc o schimbare. În schimb, clienții se înregistrează ca *subscribers* la server, astfel că serverul va anunța modificările doar clienților înregistrați pentru notificare.

Avantajul acestei abordări este că serverul nu trebuie să memoreze informații detaliate referitoare la stările clienților. Astfel, se simplifică logica aplicației și crește performanța serverului.

Pentru a minimiza traficul de rețea și a folosi tehnici de *load balancing*, unele componente ale aplicației vor fi replicate la nivel de client (de exemplu, elementele de grafică în timp real) iar altele vor fi păstrate doar pe server (de exemplu, informații de autentificare ale utilizatorilor, logica spațiului virtual, etc).

### **3.4.3 COMPONENTELE PRINCIPALE ALE APLICAȚIEI SERVER**

Aplicația server conține mai multe componente:

- Componenta de logică a spațiului virtual
  - Modulul de inteligență artificială
  - Modulul de gestiune a entităților spațiului virtual
  - **Modulul de arbitrare**
- Componenta de rețea
  - Modulul de rețea
  - Modulul de criptare/decriptare
  - Modulul de mesagerie
- Componenta de stocare a lumii virtuale
  - Modulul reprezentării lumii virtuale
  - Modulul bazei de date al lumii virtuale
- Componenta de gestiune a utilizatorilor
  - Modulul de autentificare
  - Modulul bazei de date de utilizatori

Modulul de arbitrare este componenta centrală a serverului, fiind responsabilă cu transmiterea de mesaje între diferitele module ale acestuia. Acesta gestionează ceasul principal (*tick*) al spațiului virtual 3D MMO și transmite mesaje utilizatorilor care se află în lumea virtuală. Modulul arbitru comunică direct cu componenta care gestionează lumea virtuală, cu cea care gestionează utilizatorii, cu modulele de gestiune a entităților și de inteligența artificială și cu componenta de rețea.

Componenta de stocare a lumii virtuale conține modulul reprezentării lumii virtuale și modulul bazei de date corespunzătoare. Baza de date a lumii virtuale stochează informații despre toate obiectele care sunt prezente în întreg spațiul virtual 3D. Modulul reprezentării lumii virtuale memorează un subset al informațiilor din baza de date a lumii virtuale, pentru a permite accesul rapid al modulului arbitru la acestea. Aceste informații sunt memorate într-un cache, fiind evitate astfel apelurile repetate la baza de date care sunt consumatoare de timp și afectează performanța sistemului.

Componenta de gestiune a utilizatorilor conține modulul de autentificare al utilizatorilor și baza de date a acestora. Modulul de autentificare decide pe baza

unor reguli și a unor informații din baza de date a utilizatorilor dacă un anumit utilizator se poate conecta la server. De asemenea, tot la nivelul acestei componente se stabilesc rolurile și drepturile de acces pe care utilizatorii le au în cadrul lumii virtuale. Baza de date a utilizatorilor deține informațiile utilizatorilor precum: username, parola, nume real sau informații despre personajul care reprezintă utilizatorul în cadrul lumii virtuale.

Modulul de gestiune al grupurilor spațiului virtual stochează informații despre diferitele grupuri de care poate aparține un anumit utilizator. Acest modul permite trimiterea de mesaje de către arbitru către un întreg grup de utilizatori. Utilizatorii sunt reprezentați în lista de notificări menținută în acest modul.

Modulul de inteligență artificială este responsabil cu gestiunea personajelor care nu sunt sub controlul direct al utilizatorilor, aceste personaje fiind denumite generic NPC (non player characters).

Componenta de rețea conține modulele de rețea, de mesaje și de criptare/decriptare. Modulul de mesaje se ocupă cu crearea și transmiterea mesajelor. În momentul în care modulul arbitru decide trimiterea unui mesaj către utilizatorii din lista sa de notificări, acesta este compus de către modulul de mesaje, după care este transmis modulului de criptare/decriptare. Modulul de criptare/decriptare comunică direct cu modulul rețea. De asemenea, modulul de mesaje are și rolul de a primi mesaje decriptate și de a le translata în apeluri ale modulului arbitru. Modulul de mesaje servește scopului de furnizare a unei interfețe uniforme pentru server și client și realizează totodată minimizarea numărului de pachete care sunt trimise prin rețea. Modulul rețea este responsabil cu inițierea, acceptarea sau distrugerea de conexiuni.

Clientul înglobează componentele lumii virtuale care sunt stocate pe calculatoarele utilizatorilor. Pentru a obține performanță maximă, este nevoie să se aloce clienților o cantitate cât mai mare din volumul de lucru al întregului sistem.

#### **3.4.4 COMPONENTELE PRINCIPALE ALE APLICAȚIEI CLIENT**

Aplicația client este alcătuită din următoarele componente principale :

- Componenta grafică
  - Modulul grafic
  - Modulul de vizualizare

- Modulul de gestiune a intrărilor utilizatorilor
- Componenta de rețea
  - Modulul de mesaje
  - Modulul de criptare/decriptare
  - Modulul de rețea
- Componenta de control
  - **Modulul de arbitrare**
  - Modulul de comenzi
  - Modulul de chat

Șablonul folosit în arhitectura aplicației client este model-view-controller, fiind foarte răspândit în aplicațiile care conțin interfețe cu utilizatorul. Această arhitectură nu respectă în totalitate modelul initial MVC, fiind eliminată legătura dintre model și view pentru a obține o viteză mărită.

Modulul arbitru este analog celui prezent în arhitectura serverului. Acesta este responsabil cu luarea unor decizii locale, spre deosebire de omologul său de pe server care se ocupă cu luarea unor decizii globale. Modulul arbitru este responsabil cu transmiterea mesajelor între diferitele module ale clientului și cu gestionarea ceasului local al clientului.

Componenta grafică conține modulele de grafica, de vizualizare și de gestiune a intrărilor. Modulul de vizualizare este apelat de către arbitru în vederea folosirii funcționalităților modulului grafic. Modulul grafic conține toate apelurile grafice dependente de platforma sau de API și se ocupă cu redarea 3D a spațiului virtual pe ecran. Modulul de vizualizare efectuează și translatarea intrărilor primite de la utilizator în mesaje care să poată fi înțelese de arbitru. Modulul de gestiune a intrărilor preia intrările de la utilizator și le transmite modulului de vizualizare, care determină semnificația acestora și transmite mesajele arbitrului. Acest modul abstractizează intrările provenite de la utilizator indiferent dacă acestea sunt transmise prin intermediul unei tastaturi, a unui mouse sau a altui dispozitiv de intrare.

Componenta de control conține modulele de comandă și de chat. Modulul de comandă gestionează anumite comenzi ale utilizatorului care nu au legătură directă

cu interacțiunea personajului direct cu lumea virtuală. Astfel de comenzi pot fi cele de verificare a timpului pe server și pe client. Modulul de chat este responsabil cu gestiunea comenzilor de chat prin formarea unui mesaj care este trimis arbitral, acesta urmând să îl comunice la rândul lui serverului pentru distribuție.

Componenta de rețea conține modulele de mesaje, de criptare/decriptare și de rețea. Modulul de mesaje creează mesaje care vor fi trimise serverului pentru distribuție și translatează mesajele trimise de către server clientului. Modulul de criptare/decriptare efectuează criptarea/decriptarea mesajelor trimise, respectiv primite de client. Modulul de rețea conectează aplicația client la rețeaua fizică și gestionează primirea și trimiterea de pachete.

### **3.4.5 LIMITĂRI ȘI INCONVENIENTE**

Arhitectura Client-Server distribuită pretinde că generează mai puțin trafic în rețea și evită congestionarea serverului. Totuși, au existat studii care afirmă că traficul în rețea al unui sistem centralizat este aproximativ același ca și cel în cazul unui sistem distribuit. Acest rezultat poate să fie justificat prin topologia rețelei și aplicarea diferitelor tehnici de optimizare specifice arhitecturii client-server.

O versiune la scară redusă a unui MMO poate să fie alcătuită dintr-un singur server și mai mulți clienți conectați la server. Aceasta implementare limitează numărul de jucători ce pot accesa lumea virtuală precum și dimensiunile acesteia. Pentru o versiune mai complexă și mai imersivă a spațiului virtual este necesar ca partea de server să fie implementată folosind un număr mare de mașini fizice.

Deși arhitecturile de tip Client-Server pun la dispoziție funcționalitățile necesare unei aplicații MMO, o fac însă cu un cost ridicat. Din nefericire, natura puternic centralizată a arhitecturii Client-Server introduce o sugrumare din punct de vedere al performanței. În ciuda numărului lor, care poate fi foarte mare, clusterelor de servere sunt limitate din punct de vedere computațional și traficul de rețea este concentrat practic în echipamentele din data center.

Costurile mari introduse de arhitectura Client-Server limitează practic scalabilitatea ([GUP09]) acesteia, fiind necesară împărțirea spațiului virtual în mai multe instanțe / sharduri. Cu un număr relativ mic de servere care trebuie să facă față unui volum mare de sarcini, o singură defecțiune poate să provoace căderea întregului sistem. De asemenea, pentru a obține un nivel de performanță care să facă față unui

număr mare de utilizatori, producătorii spațiilor virtuale MMO se confruntă cu un cost financiar semnificativ pentru a menține infrastructura sistemului funcțională.

## **3.5 ARHITECTURI P2P**

### **3.5.1 DESCRIERE GENERALĂ**

Aplicațiile MMO la scară redusă, de cele mai multe ori, folosesc abordări P2P. Folosirea modelului Peer-to-Peer oferă anumite avantaje într-un mod natural atrăgând astfel cercetari în domeniu în ceea ce privește această abordare pentru aplicațiile MMO ([MIL10]).

Aplicațiile la scara redusă, în mod uzual, partajează între participanți încărcarea simulării unui mediu virtual. Astfel, abordarea P2P permite utilizatorilor noi să ofere sistemului resurse suplimentare astfel încât acesta să facă față încărcării suplimentare pe care aceștia o introduc, sistemul reușind astfel să scaleze corespunzător.

Un alt avantaj al acestei abordări este că dacă una dintre resurse se defectează, ceilalți participanți pot prelua atribuțiile acesteia oferind astfel sistemului robustețe. Traficul de rețea se desfășoară în funcție de necesitățile participanților permițând partajarea încărcării între utilizatorii implicați în mod direct.

Toate aceste caracteristici oferă de asemenea reducerea costurilor pentru proprietarul spațiului virtual care practic reduce cheltuielile cu mentenanța spațiului virtual.

Principala provocare în adaptarea unei lumi virtuale 3D MMO la o arhitectură peer-to-peer o constituie realizarea funcționalităților serverelor centralizate după modelul distribuit.

Astfel, se impune menținerea consistenței unui spațiu virtual partajat de un număr foarte mare de utilizatori fără intervenția unui server. Un concept important este acela al gestiunii interesului ([BEZ08]) care susține că un utilizator singular nu trebuie să aibă cunoștințe despre ceea ce se întâmplă în lumea virtuală atâta timp cât nu îl afectează și că orice avatar are o viteză de mișcare și o capacitate de detecție limitată. Viziunea utilizatorului asupra lumii virtuale se rezumă la un

anumit domeniu de interes, acesta fiind notificat numai în privința evenimentelor care se petrec în domeniul acestuia de interes.

Schemele de gestiune a interesului pot fi clasificate folosind :

- modele spațiale
- modele de regiuni de tip publish / subscribe
- modele de comunicație hibride

Modelul spațial mai este numit și modelul “aura-nimbus” unde aura semnifică granițele care limitează prezența unui obiect în spațiu, iar nimbus presupune nivelul de conștientizare reciprocă a două obiecte. Fiecare obiect ar trebui să comunice cu obiectele care se află în *nimbusul* său în vederea pregătirii interacțiunilor viitoare. Avantajul modelului spațial este acela că numai mesajele necesare sunt transmise între entități. Totuși este nevoie ca entitățile să schimbe informații privitoare la actualizarea pozițiilor pentru a determina dacă domeniile acestora de interes se intersectează.

Determinarea de către un utilizator a vecinilor săi se realizează prin menținerea de către acesta a unei diagrame interne, construită pe baza poziției celorlalți utilizatori. Deși se reduce numărul mesajelor schimbate, se crește efortul computațional impus de menținerea diagramei interne.

Modelul regiunilor bazate pe tehnica publish / subscribe presupune împărțirea spațiului virtual în regiuni statice. Responsabilitatea mecanismului de gestiune a interesului este aceea de a determina regiunile care intersectează domeniul de interes al utilizatorului și de a forma o regiune de subscriere pentru evenimentele relevante pentru reuniunea regiunilor intersectate.

Avantajul acestui model îl constituie faptul că este mai simplu și mai ieftin de calculat regiunea de subscriere a unui utilizator, decât calculul coliziunilor domeniului de interes al acestuia. Canalul de emisie al unei regiuni se poate mapa foarte ușor peste un grup de utilizatori, astfel încât evenimentele lumii virtuale sunt transmise eficient.

Modelul de comunicare hibrid este o mixtură între modelul spațial și cel bazat pe regiuni. Lumea virtuală este împărțită în regiuni și este desemnată o super-entitate



(care rulează pe unul sau mai multe calculatoare ale utilizatorilor din regiunea respectivă) care să se ocupe de fiecare regiune. Atunci când un utilizator este pe cale să intre într-o anumită regiune, acestuia i se comunica identitatea super-entității care gestionează regiunea curentă, căreia acesta trebuie să îi transmită actualizările în ceea ce privește poziția sa.

Utilizatorii ale caror domenii de interes sunt pe cale să se intersecteze sunt notificați de către super-entitatea care gestionează regiunea în care se află, în vederea realizării unei conexiuni peer-to-peer între aceștia. Acest model profita de avantajele modelelor spațiale și bazat pe regiuni, fiind mai ușor de implementat decât modelul spațial pur.

În comparație cu arhitectura Client-Server, arhitectura Peer-to-Peer adoptă o serie de măsuri care reprezintă soluții pentru problemele întâlnite în utilizarea arhitecturii Client-Server:

- interschimbarea mesajelor se realizează direct între utilizatori, fără a mai fi nevoie de trecerea printr-un server central
- se pot folosi resursele disponibile pe mașinile utilizatorilor în ceea ce privește puterea computațională, capacitatea de stocare a informațiilor și lățimea de bandă
- descentralizarea crește robustețea sistemului deoarece se elimină punctul unic de defectare, a cărui cădere ar fi condus la încetarea funcționării întregului sistem

### **3.5.2 LIMITĂRI ȘI INCONVENIENTE**

Deși arhitectura Peer-to-Peer atunci când este folosită pentru implementarea serverelor de spații virtuale 3D MMO oferă în mod natural avantajele descrise mai sus, ridică în același timp și o serie de probleme care trebuie rezolvate :

- riscul de securitate ([WIE06]) care apare odată cu oferirea puterii de decizie aplicației ce rulează la utilizator, acesta putând modifica aplicația pentru a obține avantaje care nu îi sunt cuvenite;
- menținerea unei lumi virtuale persistente care să aibă aceleași caracteristici pentru toți utilizatorii care fac parte din ea;
- ordonarea evenimentelor și propagarea acestora către toți utilizatorii.

Atunci când se folosește pentru schema de gestiune a interesului modelul spațial, dezavantajele sunt date de dificultatea de determinare a dimensiunii optime a unei regiuni astfel încât utilizatorii să poată schimba mesaje între ei înainte să treacă în alte regiuni. Totuși, folosirea unor regiuni mari conduce la primirea de către utilizatori de mesaje nerelevante. De asemenea, o altă deficiență a modelului este aceea ca aceste regiuni nu vor funcționa bine atunci când obiectele sunt inegal distribuite.

De asemenea, folosirea modelului de comunicație hibrid poate necesita un volum mare de calcule și mesaje interschimbate pentru super-entitatea care gestionează o regiune aglomerată, fiind necesare tehnici dinamice care să distribuie sarcinile unei mulțimi de super-entități, fiecare având în gestiune o sub-regiune. O super-entitate este un potențial punct de defectare, fiind nevoie de pregătirea unor super-entități de rezervă care să preia controlul în cazul defectării celei curente.

### **3.6 ARHITECTURI HIBRIDE**

Atât arhitectura client-server cât și cea peer-to-peer prezintă unele avantaje și dezavantaje. Există unele studii ([CHE06]) care propun o arhitectură hibridă compusă din următoarele elemente:

- un server central
- mai mulți clienți selecționați ce sunt legați între ei și cu serverul într-un mod peer-to-peer
- clienți obișnuiți care se conectează folosind mecanismul client-server

În funcție de anumiți parametri ce sunt calculați folosind diverși factori (latentă, trafic generat, operații executate, etc) modul de procesare al clienților poate să fie schimbat dinamic între cele două abordări (client-server respectiv peer-to-peer).

Deși arhitectura peer-to-peer reduce semnificativ traficul de date de la serverul central, ea este vulnerabilă la posibile defectări care pot duce la inconsistente în lumea virtuală. De aceea, ar fi indicată o arhitectură hibridă în care un client devine server local atunci când are suficienta putere computațională, spațiu de stocare și lățime de bandă. Pentru a trata defectările, serverul central are o copie a stării

jocului din fiecare server de joc local.

Serverul central se ocupă doar de schimbări de stare precum mutarea într-o nouă regiune, completarea unui quest sau atacarea unui jucător. Majoritatea operațiilor sunt operații poziționale ce presupun mișcarea jucătorului prin lumea virtuală (fără trecerea în alte regiuni). Aceste operații pot fi executate de către serverele locale (clienții selecționați).

Astfel, scade traficul de date îndreptat către server, și lățimea de banda câștigată poate să fie folosită pentru mărirea numărului de utilizatori. Această arhitectură scalează mai bine decât arhitectura tradițională client-server folosind și unele din caracteristicile specifice abordării peer-to-peer.

### **3.7 GRID / CLOUD COMPUTING**

Un alt mod de abordare ce ar putea fi viabil, mai ales datorită ultimelor dezvoltări tehnologice apărute în domeniu cât și faptului ca unele cerințe, cum ar fi lățimea de bandă, devin din ce în ce mai accesibile utilizatorilor, ar fi cel care folosește pentru procesare Grid / Cloud Computing.

Astfel, partea de putere de calcul necesară arhitecturilor serverelor de spații virtuale 3D MMO ar putea fi obținută din platformele de grid / cloud computing ce sunt disponibile pe piață.

În practică însă, există provocări ce trebuie depășite pentru a se obține o soluție funcțională pentru o arhitectură de spațiu virtual 3D MMO integrată cu grid / cloud computing și din această cauză în acest moment există doar cercetări în această direcție fără a fi disponibile și soluții funcționale.

Câteva dintre provocările integrării arhitecturilor de spații virtuale 3D MMO cu soluții de grid / cloud computing ar fi :

- complexitatea modificărilor ce trebuie efectuate la nivelul arhitecturilor actuale de spații virtuale 3D MMO pentru a folosi grid / cloud computing
- latența încă destul de ridicată pentru a putea folosi în timp real grid / cloud computing pentru un volum foarte mare de operații
- toleranța la defecte

- securitatea accesului și secretul datelor folosite în cadrul spațiului virtual

La nivel de aplicații client, în ultimul an au fost lansate 2 platforme, bazate pe tehnologie de streaming, care oferă servicii de acces online la diverse jocuri clienților prin Internet : Gaikai (*[GAIII]*) și OnLive (*[ONL10]*).

Aceste platforme sunt încă în fază incipientă, numărul de aplicații puse la dispoziție este mic și sunt foarte dependente de lățimea de bandă a utilizatorilor care le accesează.

### **3.8 CONCLUZII**

Arhitectura cea mai des întâlnită folosită de aplicațiile MMO este în prezent cea Client-Server. Aceasta pune la dispoziție funcționalitățile necesare unei aplicații MMO, însă o face cu un cost ridicat. Avantajul major al unei soluții cu servere centralizate este că oferă un grad mare de control asupra sistemului, care face posibilă asigurarea autentificării utilizatorilor, persistenței și securității spațiului virtual.

Din nefericire, natura puternic centralizată a arhitecturii Client-Server introduce o sugrumare din punct de vedere al performanței. În ciuda numărului lor, care poate fi foarte mare, clusterelor de servere sunt limitate din punct de vedere computațional și traficul de rețea este concentrat practic în echipamentele din data center.

Costurile mari introduse de arhitectura Client-Server limitează practic scalabilitatea acesteia, fiind necesară împărțirea spațiului virtual în mai multe instanțe / sharduri. De asemenea, pentru a obține un nivel de performanță care să facă față unui număr mare de utilizatori, producătorii spațiilor virtuale MMO se confruntă cu un cost financiar semnificativ pentru a menține infrastructura sistemului.

Arhitecturile Peer-To-Peer, chiar dacă rezolvă unele dintre neajunsurile arhitecturilor Client-Server, au la rândul lor câteva lipsuri semnificative care le fac de nefolosit la scara largă în aplicațiile MMO moderne.

Cel mai important neajuns al arhitecturilor Peer-To-Peer este acela că nu au nici un mecanism prin care să asigure persistența stării lumii virtuale. Atunci când un utilizator se deconectează de la spațiul virtual, resursele puse la dispoziție dispar, inclusiv datele pe care aplicația acestuia le gestiona.

Faptul ca nu exista un mecanism de autoritate centralizat face foarte dificilă actualizarea spațiului virtual de către producători și introduce posibile găuri de securitate.

De asemenea, deși participanții pot pune la dispoziție resurse suplimentare sistemului, acest lucru se realizează într-un mod eterogen, fiecare mașina având propriile limite computaționale. Fără un mecanism simplu de distribuire a încărcării, o concentrare mare de activități asupra unui peer poate sa consume foarte repede resursele acestuia.

Ținând astfel cont de lipsurile și inconvenientele pe care arhitecturile actuale le prezintă , este necesară explorarea de noi soluții privind arhitecturile pentru servere de spații virtuale 3D MMO și de asemenea găsirea unor modalități de optimizare a operațiilor efectuate în cadrul simulării unui lumi virtuale pentru a elimina, sau măcar reduce, o parte din problemele curente cu care se confruntă aplicațiileMMO.

## 4 TEHNICI INOVATIVE DE PARALELISM PENTRU SPAȚII VIRTUALE 3D MMO BAZATE PE GPGPU

GPGPU (General Purpose on Graphical Processing Units) reprezintă o metodă prin care unitatea de procesare grafică poate fi utilizată pentru a executa calcule/programe de uz general. Acest lucru a devenit posibil odată cu apariția etapelor programabile ale GPU precum și datorită existenței bibliotecilor și instrumentelor de dezvoltare puse la dispoziție de producătorii de unități de procesare grafică.

### 4.1 DE CE GPGPU

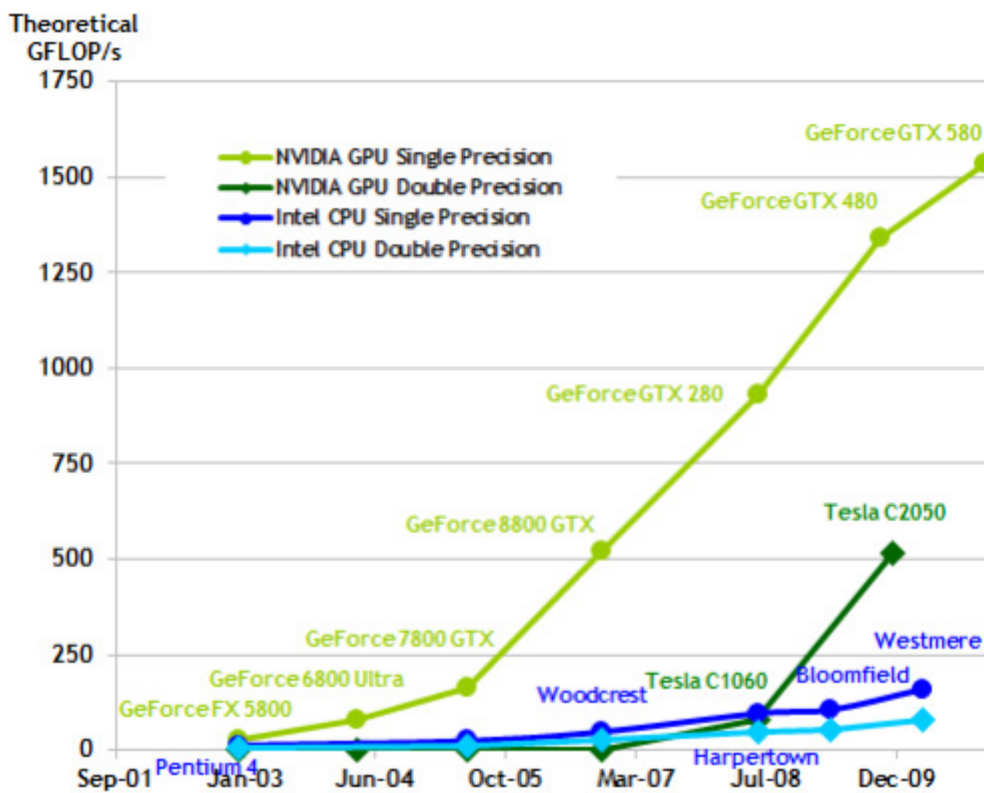
Unitatea de Prelucrare Grafică (GPU) modernă a devenit un hardware foarte versatil în domeniul arhitecturilor de calcul paralele multi-core ([BRE08], [CHE08]). Aceste arhitecturi, care includ GPU și CPU multi-core de la AMD și Intel, procesoarele CELL, procesoarele SUN UltraSparc, se diferențiază față de arhitectura clasică CPU prin următorul fapt: favorizarea operațiilor ce se pot executa în paralel asupra unei cantități mari de date față de execuția operațiilor *single-task* de latență scăzută.

Type	Processor	Cores/Chip	ALUs/Core
GPU	AMD Radeon HD 5870	20	80
	NVIDIA GeForce GTX 480	15	32
CPU	Intel XEON Westmere	6	2
	CELL	8	4

*Figura 4-1: Arhitecturi GPU și CPU moderne*

Deși există diferențe de implementare între diverșii producători, toate GPU-rile moderne încearcă să mențină o eficiență ridicată prin folosirea de arhitecturi multi-core, care folosesc atât multithreading hardware cât și procesare SIMD – Single Instruction Multiple Data. Aceste tehnici nu sunt unice pentru GPU însă în comparație cu CPU, GPU duc la extrem aceste arhitecturi.

De exemplu, GPU NVIDIA GeForce GTX590 are 1024 de procesoare scalare care operează la 1,2 GHz. Acestea sunt organizate în grupuri de 32 și au un *peak-rate* (vârf de performanță) de 2,5TFLOPS. În comparație, un procesor high-end de la Intel , XEON Westmere care operează la o frecvență de 3.3GHz conține 6 core-uri și are un *peak-rate* de 146GFLOPS.



*Figura 4-2: Comparație între NVidia GPU și Intel CPU*

Sursa : NVidia CUDA SDK 4.0

### 4.1.1 ABSTRACTIZĂRI CUDA

CUDA (Compute Unified Device Architecture) este o arhitectură hardware și software pentru realizarea și gestionarea calculelor pe GPU, fără a fi nevoie de maparea la un API specializat pentru grafică.

În programarea realizată cu ajutorul CUDA, GPU este văzut ca un co-procesor la CPU principal, numit gazdă. Cu alte cuvinte, porțiuni care sunt intensiv computaționale și prezintă paralelism de date, din aplicația ce rulează pe gazdă, sunt încărcate pe dispozitiv (GPU).

Această porțiune din aplicație este organizată sub forma unei funcții care este executată simultan de un număr mare de fire de execuție. Funcția destinată rulării pe GPU este compilată în instrucțiuni specifice dispozitivului GPU, rezultând un program numit *kernel*.

CUDA furnizează 3 abstracții principale:

- ierarhie de grupuri de fire de execuție,
- spații de memorie
- sincronizarea prin bariere

Aceste abstracții ușurează comunicarea și sincronizarea între firele de execuție care reprezintă principalele unități computaționale. Astfel, programatorul este îndrumat să își partiționeze problema în subprobleme de granularitate mare, care la rândul lor sunt divizate în subprobleme de granularitate mai mică ce pot fi rezolvate **independent** și **cooperativ**.

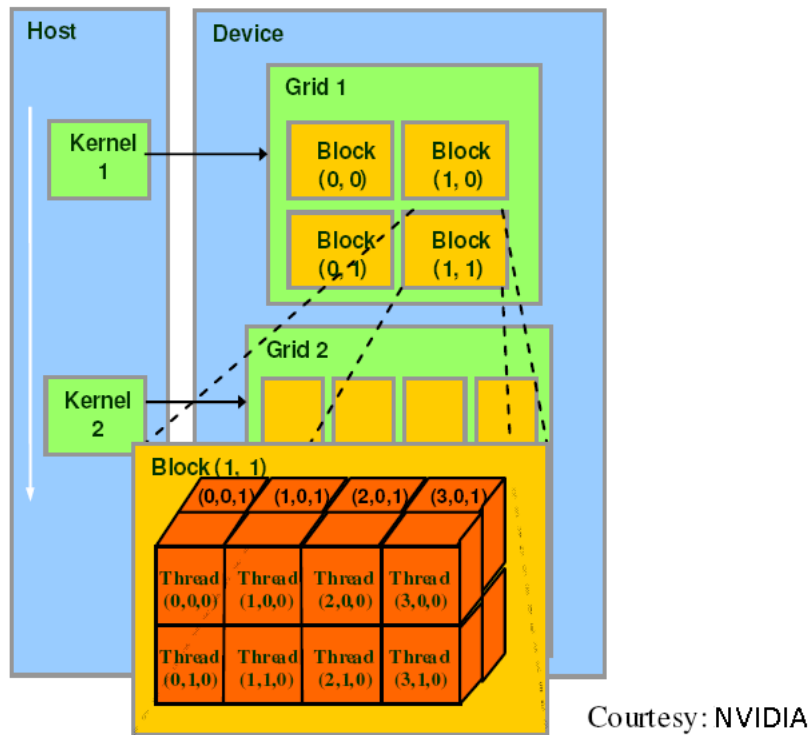
#### 1. Ierarhie de grupuri de fire de execuție.

Un grup de fire de execuție care cooperează prin partajarea eficientă a datelor și își pot sincroniza execuția poartă numele **de bloc de fire de execuție**. Fiecare fir de execuție dintr-un bloc este identificat printr-un număr unic la nivelul blocului. Un bloc mai poate fi definit și ca un tablou de fire de execuție cu 2 sau 3 dimensiuni. Dimensiunea unui bloc este limitată, însă se pot defini mai multe blocuri de aceeași dimensiune, care sunt **grupate într-un grid**.

Deoarece toate firele de execuție dintr-un grid executa aceeași funcție kernel, pentru a se deosebi între ele și a identifica porțiunile datelor pe care trebuie să le prelucreze, se folosește un mecanism de coordonate unic, ce permite identificarea



exacta a fiecărui fir de execuție în cadrul gridului ( **blockId** ) și în cadrul blocului ( **threadId** ).



*Figura 4-3: Ierarhie grupuri fire de execuție*

Firele de execuție din blocuri diferite nu pot însă comunica între ele și nu se pot sincroniza. Dimensiunea și forma blocurilor și grid-ului fac parte din configurația de execuție a kernelului și nu se pot modifica pe parcursul execuției.

Blocurile de fire de execuție trebuie să se execute independent. Este necesar să existe posibilitatea ca ele să se execute în orice ordine, în paralel sau în serie. Independența permite blocurilor de fire de execuție să fie planificate în orice ordine și executate pe oricâte procesoare, permițând programatorului să scrie cod scalabil. Numărul de blocuri de fire de execuție dintr-un grid este dictat în mod obișnuit de dimensiunea datelor care trebuie procesate și nu de numărul de procesoare din sistem, număr care poate fi cu mult depășit.

Un grid de blocuri este executat prin rularea unuia sau a mai multor blocuri pe fiecare multiprocesor folosind *time-slicing*. Fiecare bloc este împărțit în grupuri SIMD de fire de execuție numite *warp-uri*, fiecare având același număr de fire de

execuție, numit *dimensiunea warp-ului*. Un warp este executat pe un multiprocesor după modelul SIMD. Un planificator de fire de execuție comută de la un warp la altul pentru a maximiza folosirea resurselor computaționale ale multiprocesorului.

Modul în care un bloc este împărțit în warp-uri este mereu același: fiecare warp conține fire de execuție consecutive cu threadId-uri crescătoare, primul warp începând cu firul de execuție cu threadId-ul 0. Un bloc este procesat de un singur multiprocesor, de aceea spațiul de memorie partajată rezidă în memoria shared on-chip a multiprocesorului, rezultând astfel viteze mari de acces.

În cazul în care numărul de fire de execuție din bloc este mai mare decât numărul de registre ale multiprocesorului, kernelul nu va putea fi executat și nu va fi lansat. Mai multe blocuri pot fi procesate de același multiprocesor concurent prin alocarea registrelor și a memoriei partajate între blocuri

## 2. Spații de memorie

A două abstractizare majoră folosită de CUDA se referă la diferitele tipuri de memorie și gradul lor de accesibilitate. Ele diferă prin strategia de acces precum și prin vizibilitatea în cadrul ierarhiei de fire de execuție. Următoarele tipuri de memorie sunt vizibile pentru fiecare fir de execuție la nivel de grid :

- **registri :**
  - cea mai rapidă memorie;
  - accesibilă doar în interiorul firului de execuție; același *scope* ca firul de execuție.
- **memoria partajată :**
  - împărțită în module de memorie de dimensiune egală numite *bancuri de memorie*, fiecare banc memorând o variabilă de 32 de biți; conflictele la nivel de bloc apar atunci când sunt accesate date din cadrul aceluiași bloc, caz în care hardware-ul serializează accesul la date forțând firele de execuție să aștepte până când cererile de acces la memorie au fost rezolvate;

- în cazul în care toate firele de execuție citesc de la aceeași adresă de memorie partajată, se evită serializarea prin utilizarea unui mecanism de broadcast;
  - la fel de rapidă ca regiștrii atunci când nu exista conflicte între bancurile memoriei partajate;
  - vizibila de către toate firele de execuție din interiorul blocului; același scope ca al blocului.
- **memoria globală :**
- memoria globală și memoria locală nu beneficiază de mecanism de caching, deci fiecare acces la memoria globală sau locală generează un acces explicit la memorie; un acces la memoria locală sau la memoria globală impune o latență suplimentară de 400-600 cicli de ceas și astfel memoria globală este de aproximativ 150 de ori mai lentă decât memoria partajată;
  - dispozitivul este capabil să citească 32, 64 sau 128 de biți într-o singură instrucțiune însă datele trebuie să fie aliniat corespunzător; altfel, compilatorul va genera mai multe instrucțiuni de load;
  - operațiile de accesare a memoriei globale, din interiorul firelor de execuție ale unui warp, ar trebui aranjate astfel încât să fie reunite într-o singură instrucțiune load;
  - accesibila atât din gazdă cât și din dispozitiv; are scope-ul aplicației.
- **memoria locală :**
- nu dispune de un dispozitiv hardware dedicat, fiind alocată de către compilator în cadrul memoriei globale, având aceleași performanțe ca aceasta; variabilele unui fir de execuție pot fi plasate în memoria locală atunci când nu există un număr mare de variabile registru sau când datele ar ocupa prea mult spațiu de memorie dacă ar fi plasate în regiștri;
  - considerabil mai lentă decât regiștrii și memoria partajată;

- accesibilă doar în interiorul firului de execuție; același scope ca firul de execuție;
- spațiile de memorie locală și globală sunt regiuni read-write din memoria dispozitivului și nu dispun de cache.

- **memoria constantă**

- poate fi doar citita din interiorul unui kernel și este optimizată pentru cazul în care toate firele de execuție citesc din aceeași locație de memorie;
- memoria constantă poate fi scrisă doar de către gazdă prin intermediul funcției `cudaMemcpyToSymbol` și este persistentă pe durata execuției kernelurilor din cadrul aceleiași aplicații; se pot memora 64 de KB de date în cache-ul constant și există 8 KB de cache pentru fiecare multiprocesor;
- durata accesului la date variază de la un ciclu de ceas în cazul unui cache hit până la câteva sute de cicluri de ceas, depinzând de localitatea memoriei.

- **memoria de texturi**

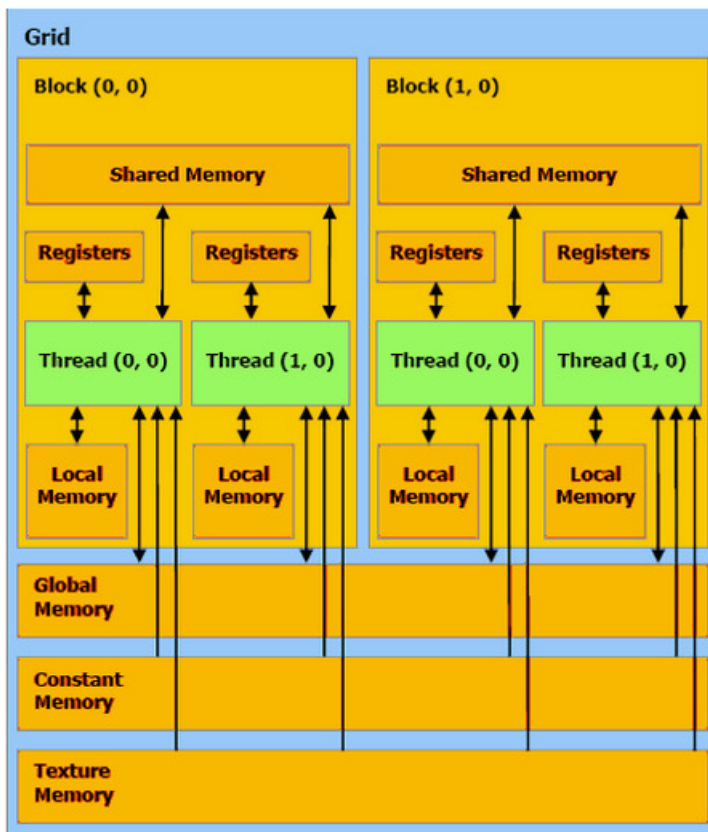
- oferă capabilități de caching a memoriei globale;
- reprezintă o modalitate de interacțiune cu capabilitățile de afișare ale GPU, având rolul de a accelera procesul de mapare și deformare a unei texturi 2D pe un model poligonal 3D; de asemenea, reprezintă o modalitate mai simplă de accesare a memoriei globale, deoarece fiecare unitate de texturi are o memorie internă care stochează date din memoria globală;
- utilizarea memoriei de texturi oferă o performanță sporită în cazul în care firele de execuție ale unui *warp* citesc de la adrese de memorie apropiate; o scădere de performanță poate apărea doar atunci când se înregistrează un *cache miss* care impune citirea de date din memoria globală;
- utilizarea memoriei de texturi poate avea următoarele beneficii :

- datele împachetate pot fi difuzate în variabile diferite într-o singură operație
- întregii pe 8 și 16 biți pot fi convertiți în numere în virgula mobilă pe 32 de biți cu valori între 0.0 și 1.0 sau -1.0 și 1.0 de către unitățile de texturare prin operații de interpolare liniară, bi-liniară și tri-liniară care sunt realizate direct de către hardware dedicat
- utilizarea memoriei de texturi pe post de cache este ușor de folosit deoarece aceasta este optimizată pentru o localitate spațială 2D și poate oferi o creștere de performanță când firele de execuție ale unui warp accesează zone aflate în locații apropiate din interiorul texturii; pentru a putea folosi memoria de texturi este nevoie să se mapeze o textură la porțiunea de memorie care se dorește a fi folosită;
- tipuri de memorie cu care poate fi asociată o textură :
  - memorie liniară : alocată cu `cudaMalloc` ; permite scrierea în memoria globală
  - tablouri CUDA : alocate cu `cudaMallocArray` sau `cudaMalloc3D`; nu este permisă scrierea în tablouri din interiorul kernelului
  - memorie 2D pitch linear : alocată cu `cudaMallocPitch`; permite scrierea în memoria globală;
- pașii care trebuie urmați pentru folosirea memoriei de texturi în CUDA sunt :
  - pe gazdă (CPU) :
    - alocare/obținere memorie (memorie liniară, pitch linear sau tablou CUDA)
    - creare obiect referință la textură
    - mapare referință textură la memorie
    - eliberare referință textură și resurse

- pe dispozitiv (GPU) :
  - *fetch* utilizând referința la textură
  - se folosește *tex1Dfetch* pentru memorie liniară și *tex1D*, *tex2D*, *tex3D* pentru tablouri și *memorie pitch linear*

Se pot trage următoarele concluzii referitoare la utilizarea diverselor tipuri de memorie în cadrul arhitecturii CUDA :

- spațiul de memorie local și global sunt implementate ca regiuni read-write ale memoriei dispozitivului și nu beneficiază de caching, ceea ce înseamnă că fiecare acces la memoria globală sau fizică generează un acces fizic
- pentru un multiprocesor, o instrucțiune de acces la memorie în cadrul unui warp durează 4 cicluri de ceas, în vreme ce un acces la memoria locală sau globală are o latență de 400-600 de cicluri de ceas
- cuvântul cheie *\_\_device\_\_* specifică o variabilă care se află în memoria globală și care nu poate fi accesată de către codul rulat pe gazdă; pentru a mari viteza de execuție se încearcă evitarea accesurilor la memoria locală sau globală și se încurajează utilizarea la maximum a variabilelor registru, *\_\_shared\_\_* sau *\_\_constant\_\_*
- pentru accesarea memoriei de texturi fiecare multiprocesor folosește o unitate de textură



*Figura 4-4: Spații de memorie*

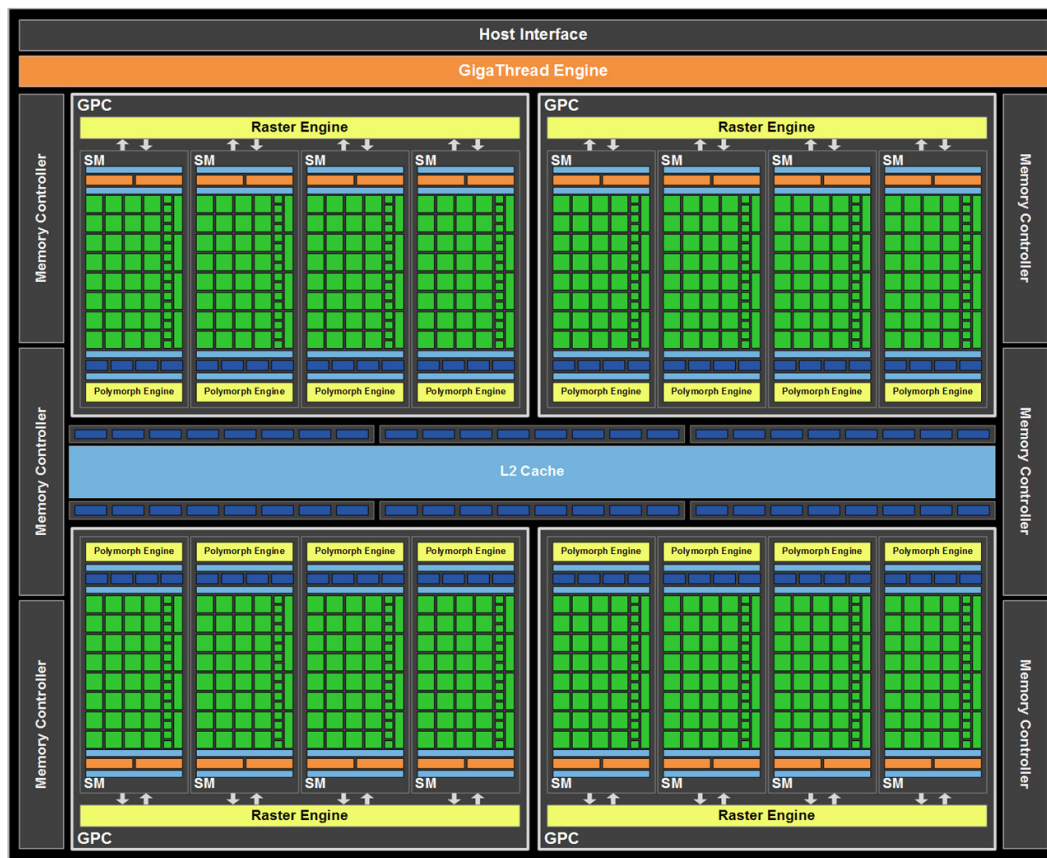
### 3. Sincronizarea prin bariere

CUDA pune la dispoziție un mecanism special de sincronizare și anume prin plasarea unor bariere. Astfel, toate firele de execuție dintr-un bloc vor fi blocate până când ultimul dintre ele atinge bariera. Costul asociat este foarte scăzut, nefiind întârzieri semnificative la apelarea instrucțiunii de sincronizare. Sincronizarea firelor de execuție în afara unui bloc nu este posibilă.

#### 4.1.2 MODELUL DE EXECUȚIE CUDA

Arhitectura procesoarelor grafice care suportă CUDA constă într-un set de multiprocesoare cu memorie partajată on-chip. Aceste multiprocesoare suporta mai multe fire de execuție și poartă numele de Streaming Multiprocessors (SM). Când un program gazdă invocă un grid de kernel-uri, blocurile grid-ului sunt distribuite către multiprocesoare. Firele de execuție dintr-un bloc se execută concurent pe un

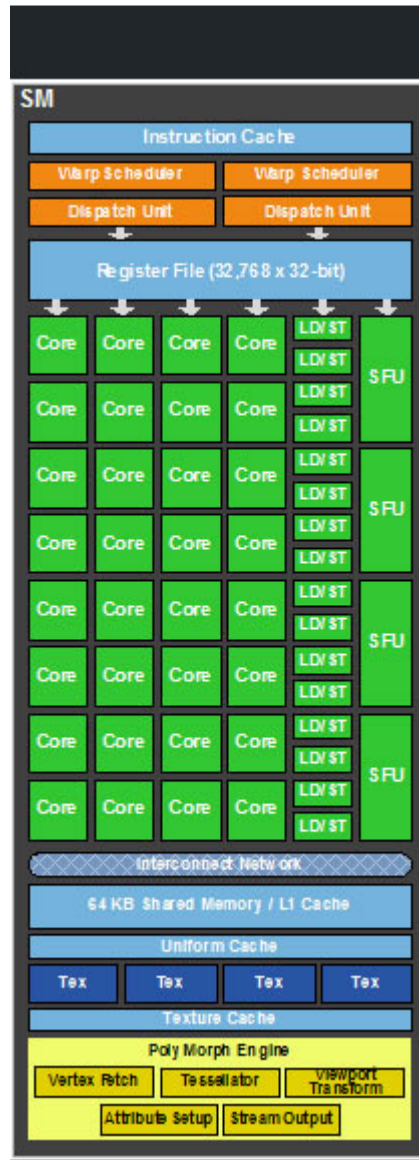
multiprocesor. Când firele de execuție din blocuri se termină, noi blocuri sunt lansate pe multiprocesoarele rămase libere. Chipsetul 580GTX de la NVidia, care este construit pe cea mai nouă arhitectură de GPU-uri denumită Fermi, conține 16 astfel de multiprocesoare.



*Figura 4-5: Arhitectura Fermi*

Un multiprocesor conține 32 de procesoare scalare (SP), două unități speciale pentru funcții, o unitate de instrucțiuni cu mai multe fire de execuție, și memorie partajată on-chip. Multiprocesoarele creează, gestionează și execută firele de execuție concurente în hardware fără penalizări date de planificare. Bariera pentru sincronizare este implementată printr-o singură instrucțiune cea ce face sincronizarea foarte rapidă. Astfel se permite paralelism de granularitate fină, descompunându-se problema prin atribuirea unui fir de execuție fiecărui element de date.





**Figura 4-6: Conținut SM din arhitectura Fermi**

Multiprocesorul mapează fiecare fir de execuție la un procesor scalar, astfel execuția se va face independent folosindu-se adrese de instrucțiuni și registre separate. Două sau mai multe blocuri pot rula pe același multiprocesor. Dacă un bloc de fire de execuție este blocat în așteptarea accesului la memoria globală, de exemplu, un alt bloc poate fi rulat.

Numărul de blocuri pe care un multiprocesor le poate procesa la un moment dat depinde de câte registre pe fir de execuție și câtă memorie partajată pe bloc sunt

necesare pentru un kernel, deoarece registrele multiprocesorului și memoria partajată sunt resurse folosite de către toate firele de execuție ale blocurilor active. Pentru gestionarea numărului mare de fire de execuție multiprocesoarele folosesc o nouă arhitectură numită SIMT (single-instruction, multiple-thread).

Unitatea SIMT din multiprocesor, creează, gestionează, planifică și execută firele de execuție în grupuri de câte 32, numite warps. Firele de execuție individuale ce compun un warp, pornesc împreună la aceeași adresă de program, dar sunt libere să se ramifice și să se execute independent.

La fiecare moment la care se emite o instrucțiune, unitatea SIMT selectează un warp care este gata să se execute și emite următoarea instrucțiune către firele de execuție active ale warp-ului. Un warp execută o instrucțiune comună la fiecare moment de timp, deci eficiența maximă este atinsă când toate cele 32 de fire de execuție cad de acord asupra căii lor de execuție.

Dacă firele de execuție dintr-un warp diverg din cauza unor ramuri condiționale dependente de date, warp-ul execută serial fiecare cale luată, dezactivând firele de execuție care nu sunt pe acea cale și când toate căile se termină, firele de execuție converg înapoi pe aceeași cale de execuție. Divergența ramurilor are loc doar în cadrul unui warp. Warp-urile diferite se execută independent, indiferent dacă execută căi comune sau disjuncte.

### 4.1.3 ELEMENTE DE LIMBAJ CUDA

Programarea în CUDA presupune izolarea unei porțiuni dintr-o aplicație care se execută de multe ori pe date diferite, în mod independent, într-o funcție care va fi rulată pe GPU. Aceasta funcție poartă denumirea de *kernel* și este compilată în setul de instrucțiuni specific device-ului, după care este transmisă acestuia spre execuție.

O funcție kernel este o funcție care este apelată de pe gazdă și este executată pe GPU de către mai multe fire de execuție simultan.

Apelul unei funcții kernel conține numele acesteia, configurația de execuție separată de “<<<” și “>>>” și parametrii funcției între paranteze rotunde.

Configurarea de execuție conține numărul de blocuri și numărul de fire de execuție din fiecare bloc.

La definiția funcției kernel se prefixează antetul acesteia cu cuvântul cheie `__global__` care specifică o funcție kernel care se execută pe GPU, dar care poate fi apelată numai de pe sistemul gazdă. De asemenea, o funcție kernel trebuie să întoarcă întotdeauna valori de tipul `void`.

Atât gazda cât și dispozitivul au propria memorie DRAM numită memoria gazdei, respectiv memoria dispozitivului. Pentru a putea prelucra datele pe GPU este necesar transferul acestora de pe CPU pe GPU, acest lucru fiind realizat prin intermediul **cudaMemcpy**. Însa, înainte de a transfera datele pe GPU este nevoie să se aloce memorie pentru acestea folosind **cudaMalloc**.

De asemenea, transferul rezultatelor pe CPU se realizează tot prin intermediul `cudaMemcpy`, tipul de transfer al datelor (`cudaMemcpyHostToHost`, `cudaMemcpyHostToDevice`, `cudaMemcpyDeviceToHost`, `cudaMemcpyDeviceToDevice`) fiind specificat printr-un parametru.

CUDA oferă posibilitatea folosirii limbajelor de programare de nivel înalt cum ar fi C/C++ pentru a dezvolta aplicații care beneficiază de avantajul unui nivel ridicat de performanță și scalabilitate pe care arhitecturile GPU le oferă. Scopul extensiei limbajului C oferită de CUDA este de a furniza acomodarea ușoară pentru utilizatorii familiari cu limbajul de programare C. Există și unele limitări pentru codul ce va fi rulat pe dispozitiv, nefiind suportată recursivitatea și nici pointerii la funcții, dar există suport pentru șabloane, referințe și supraîncărcarea funcțiilor și a operatorilor asemănătoare cu cele din C++.

Extensiile de limbaj se pot grupa în patru categorii:

- Atribute pentru tipuri de funcții pentru a se specifica dacă funcția se execută pe gazdă sau pe dispozitiv, și dacă ea poate fi apelată de către gazdă sau numai de către dispozitiv. Exemple: `__device__`, `__global__`, `__host__`.
- Atribute pentru tipul variabilelor pentru a specifica locația de memorie unde va fi stocată o variabilă. Exemple: `__device__` (spațiul de memorie global), `__constant__` (spațiul de memorie constant), `__shared__` (spațiul de memorie partajată).

- Variabile predefinite care specifică dimensiunile blocului și grid-ului precum și indicii firelor de execuție și ai blocurilor. Exemple: `blockIdx`, `gridDim`, `threadIdx`, `blockDim`, `warpSize`.

Exemplu pentru calculul identificatorului firului de execuție în cadrul grid-ului și blocului :

```
__global__ void actiune_fir de execuție(float* a_d, float* b_d,float *r_d,int N)
{
    // Identificarea poziției exacte a firului de execuție
    int x,y;
    x = (blockIdx.x * blockDim.x + threadIdx.x);
    y = (blockIdx.y * blockDim.y + threadIdx.y);
    ...
}
```

Folosind atributul `__global__` funcția de mai sus se va executa pe dispozitiv și va fi apelată din programul gazdă. Pentru calcularea identificatorului se folosesc variabilele predefinite `blockDim`, `blockIdx` și `threadIdx`. Fiecare fișier sursă care conține aceste extensii trebuie compilat cu un compilator special oferit de CUDA și anume `nvcc`. Modul de operare al `nvcc` constă în separarea codului pentru dispozitiv, de codul rulat pe gazdă și compilarea codului pentru dispozitiv într-o formă de limbaj de asamblare (cod `ptx`) sau în formă binară (obiect cubin).

Aplicațiile pot să ignore codul generat pentru gazdă și să încarce și să execute codul `ptx` sau obiectul cubin pe dispozitiv, folosind CUDA driver API. Pentru codul rulat pe dispozitiv, compilatorul nu suportă opțiunii specifice C++ cum ar fi clasele, moștenirea, sau declararea variabilelor în blocuri de bază.

### 4.1.4 INTEGRARE CUDA CU BUFFEROBJECTS OPENGL

#### 4.1.4.1 PBO (Pixel Buffer Object)

CUDA se poate folosi în combinație cu OpenGL prin utilizarea PBO (Pixel Buffer Object) pentru a manipula bufferul de imagine. Procesările sunt realizate în cadrul CUDA pixel cu pixel și sunt redade pe ecran cu ajutorul OpenGL.

De asemenea, CUDA poate genera mesh-uri 3D care pot fi redată prin utilizarea VBO-urilor (Vertex Buffer Objects) OpenGL sub forma de suprafețe colorate, imagini wireframe sau mulțimi de puncte 3D.

CUDA mapează buffere-le OpenGL într-un spațiu de memorie propriu prin intermediul funcției **cudaGLMapBufferObject**. OpenGL nu trebuie să opereze modificări asupra niciunui buffer care este mapat în spațiul de memorie CUDA. Buffer-ul poate fi accesat prin intermediul unui pointer și se pot efectua modificări din interiorul kernelului, lucrul cu acesta realizându-se la nivel de pixel.

Pentru a putea mapa buffer-ul OpenGL în spațiul de memorie CUDA și a putea efectua modificări asupra acestuia prin intermediul kernelului trebuie parcurși următorii pași:

- se creează o fereastră
- se creează un context OpenGL
- se setează viewport-ul și sistemul de coordonate OpenGL
- se generează unul sau mai multe buffere OpenGL care vor fi partajate cu CUDA
- se înregistrează buffer-ele în CUDA

Crearea unui buffer se realizează prin intermediul funcției **glGenBuffers**, iar legarea acestora se realizează prin **glBindBuffer**. Alocarea memoriei se face prin apelul **glBufferData**, iar înregistrarea în CUDA se face prin apelul funcției **cudaGLRegisterBufferObject**.

Pentru a manipula bufferul de imagine dintr-o aplicație CUDA este necesară parcurgerea următorilor pași:

- se alocă un buffer OpenGL cu dimensiunea egală cu a bufferului de imagine
- se alocă o textură OpenGL cu dimensiunea egală cu a bufferului de imagine
- se mapează buffer-ul OpenGL în memoria CUDA
- se scrie imaginea din CUDA în buffer-ul OpenGL mapat
- se demapează buffer-ul OpenGL
- se atașează textura buffer-ului OpenGL
- se desenează un quad care specifică coordonatele texturii în fiecare colț
- se interschimbă bufferele pentru a se desena imaginea pe ecran

#### **4.1.4.2 VBO (Vertex Buffer Object)**

VBO (Vertex Buffer Object) se folosește pentru manipularea de obiecte 3D în CUDA și redarea acestora folosind OpenGL sub forma unei mulțimi de puncte, imagini wireframe sau suprafețe.

Pentru a crea și afișa obiectele 3D folosind VBO-uri, sunt necesare tablouri diferite pentru reținerea vârfurilor și culorilor. Utilizarea VBO-urilor furnizează performanțe mai bune atunci când buffer-ele OpenGL sunt mapate în spațiul de memorie CUDA.

Pentru utilizarea VBO trebuie parcurși următorii pași:

- crearea unei ferestre OpenGL
- setare coordonate de vizualizare
- selectare dispozitiv CUDA
- creare buffere VBO pentru vârfuri și culori
- înregistrare buffere VBO

#### **4.1.5 EVOLUȚIE CAPABILITĂȚI CUDA**

Primele versiuni de CUDA aveau următoarele caracteristici:

- integrare CUDA în driver-ul GPU
- API asincron pentru copieri de memorie și lansări de kerneluri în execuție
- API pentru interogarea statusului apelurilor CUDA
- suport pentru depanarea programelor folosind gdb
- suport pentru Visual Studio 2008
- compilare Just-in-time (JIT) pentru aplicații care generează în mod dinamic kerneluri CUDA
- suport pentru template-uri C++ în kernelurile CUDA

Toolkit-ul 2.2 oferă față de versiunile anterioare următoarele caracteristici:

- Visual Profiler pentru GPU : CUDA Visual Profiler este un instrument grafic care activează analiza performanțelor kernelurilor care rulează pe GPU. Acesta dispune de instrumente de analiză pentru tranzacțiile cu

memoria, oferind dezvoltatorilor posibilitatea de a analiza timpul consumat pentru acest tip de operații

- îmbunătățirea interoperabilității cu OpenGL - performanță îmbunătățită pentru aplicațiile OpenGL care sunt executate pe GPU-uri profesionale din seria Quadro prin implementarea calculelor folosind CUDA
- zero-copy – permite funcțiilor CUDA să scrie și să citească direct în/din memoria sistemului, reducându-se dimensiunea și frecvența cu care datele sunt transmise între CPU și GPU
- Pinned Shared System – aplicațiile multi-GPU pot avea o performanță crescută datorită posibilității de a accesa aceleași date din memoria sistemului de către mai multe GPU-uri
- memcpy asincron pe Vista – permite copierea datelor asincron, lucru care era disponibil pe alte platforme, dar nu și pe Vista
- debugger hardware pentru GPU – CUDA GDB debugger pentru Linux oferă toate caracteristicile necesare pentru a depana direct pe GPU, inclusiv setarea de breakpointuri, inspectarea de variabile sau a stării.
- modul device exclusiv - este prezentă o opțiune de folosire exclusivă a GPU care garantează ca toată puterea de procesare și memoria GPU este alocată aplicației respective
- schimbă paradigma în care datele erau copiate pe GPU, erau prelucrate pe GPU după care erau transferate pe host cu o nouă paradigmă în care se permite maparea memoriei gazdei în memoria device-ului prin intermediul funcției `cudaMallocHost`

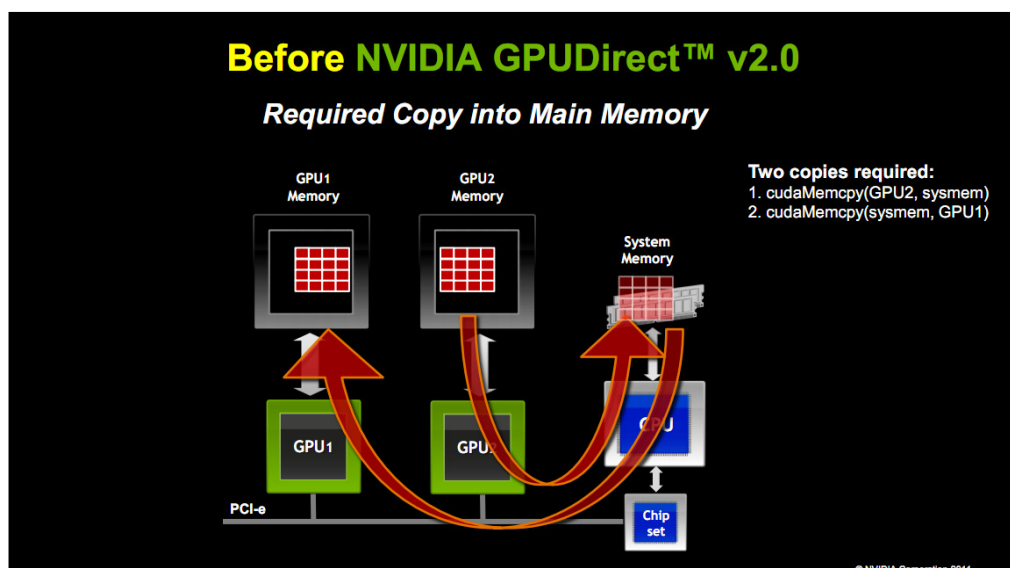
CUDA 3.0 aduce ca elemente de noutate față de versiunile anterioare următoarele caracteristici:

- interoperabilitate ridicată între CUDA și API-urile 3D OpenGL și DirectX
- CUDA-GDB dispune de un dispozitiv de verificare a memoriei în vederea adreselor nealiniat
- moștenirea claselor și template-urilor C++
- sporirea capabilităților OpenCL
- **suport pentru arhitectura Fermi**
- permite atât programarea folosind API-ul de runtime cât și cel de driver, care este mai greu de programat și presupune scrierea unei cantități mai mare de cod, dar care oferă un nivel mai bun de control și este independent de limbaj

- bibliotecile toolkit-ului CUDA suporta versionarea, foarte importantă pentru dezvoltarea de cod care să ruleze pe o varietate de sisteme fara a fi necesara recompilarea; acest lucru se realizează prin inspectarea versiunii driver-ului și oferirea caracteristicilor necesare în funcție de acestea.
- conține o versiune separată pentru CUDA C runtime pentru debugging în modul de emulare
- oferă o performanță îmbunătățită pentru programele compilate în modul de debugging deoarece registrele nu mai sunt forțate a fi reținute în memoria locală

CUDA 4.0 aduce ca elemente de noutate față de versiunile anterioare următoarele caracteristici:

- suport mult mai bun pentru aplicații multi-GPGPU :
  - o posibilitatea de partaja accesul la GPU-uri din contextul mai multor fire de execuție independente
  - o utilizarea tuturor GPU-urilor din sistem în mod concurent din cadrul unui singur fir de execuție
  - o **acces Peer-To-Peer** : posibilitatea de comunicare directă între GPU-urile din sistem



*Figura 4-7: Comunicare inter-GPU înainte de CUDA 4.0.*

Sursa: NVidia 2011



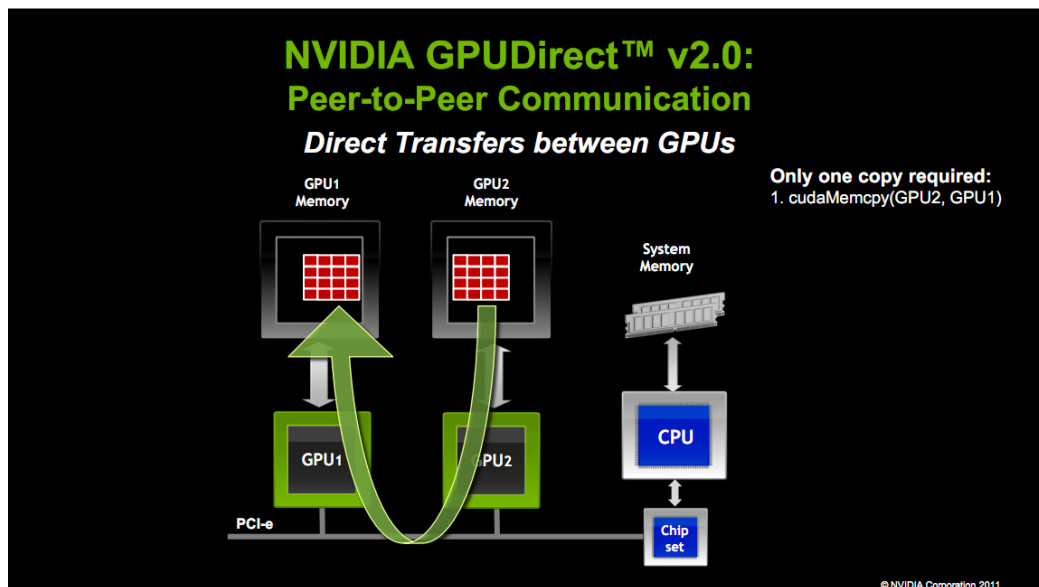


Figura 4-8: Comunicație inter-GPU Peer-To-Peer în CUDA 4.0.

Sursa: NVidia 2011

- **spațiu de memorie virtual unificat** : existență unui singur spațiu de adresare (CPU și GPU) pentru o utilizare mai facilă a resurselor de memorie

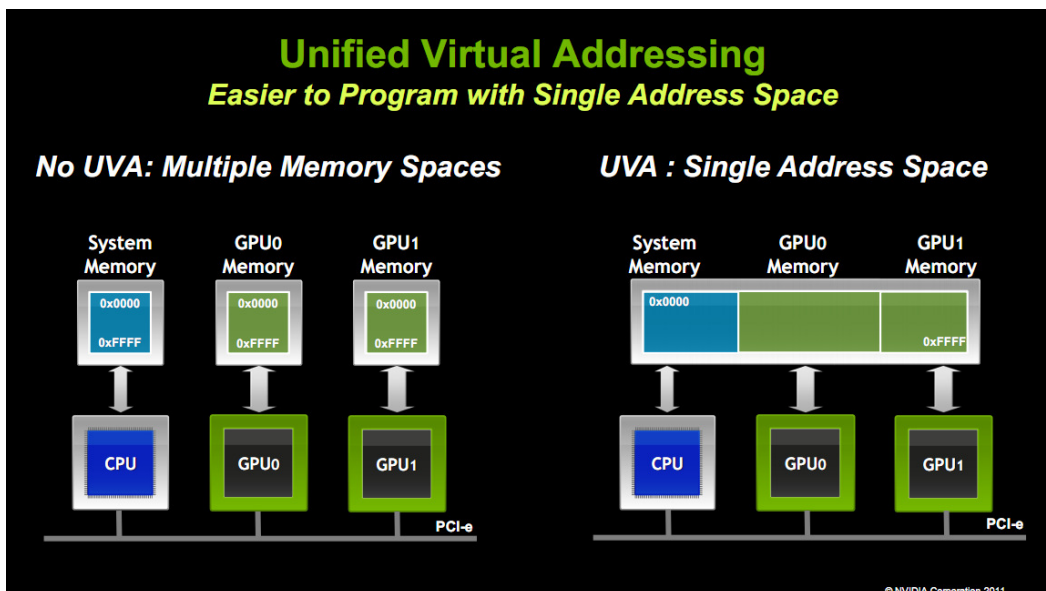


Figura 4-9: Spațiu de adresare unificat în CUDA 4.0.

Sursa: NVidia 2011

- no-copy pinning pentru memoria sistem : o alternativă mai rapidă pentru cudaMallocHost
- suport pentru operatorii C++ new și delete
- suport pentru funcții virtuale
- punerea la dispoziție a unei biblioteci cu funcții dedicate pentru procesarea de imagini

## **4.2 UTILIZAREA GPGPU PENTRU REDAREA 3D CLIENT-SIDE**

Așa cum s-a menționat și la începutul lucrării, unul dintre obiectivele principale ale unei aplicații MMO 3D pe partea de client, este acela de a crea un grad de imersiune cât mai ridicat pentru utilizatorii spațiului virtual. Aceste aplicații utilizează pe partea de client un motor grafic 3D modern care folosește multe din tehnologiile prezentate anterior pentru a crea un mediu virtual extrem de interactiv și realist.

RayTracing reprezintă o metodă de redare, care utilizată în cadrul simulărilor 3D, poate să ofere un grad de realism net superior comparativ cu tehnologiile client-side prezentate anterior care utilizează ca metodă de redare banda grafică.

Până de curând, utilizarea RayTracing ca metodă de redare în timp real nu era posibilă deoarece necesita un efort computațional ridicat și hardware-ul existent nu era suficient de rapid pentru a obține rezultate în timp real.

În continuare se va prezenta o descriere a algoritmului RayTracing, o soluție pentru a executa algoritmul RayTracing folosind procesare GPU și multi-GPGPU precum și evaluarea rezultatelor obținute.

### **4.2.1 SIMULAREA FOTO-REALISTĂ ÎN TIMP REAL A SCENELOR 3D FOLOSIND RAYTRACING**

Comparativ cu metodele de redare clasice, prin utilizarea algoritmului de RayTracing se obțin, în mod natural, efecte dificil de implementat folosind metodele clasice :

- umbre

- reflexii
- refracții

Toate aceste efecte sunt însă obținute cu un cost computațional ridicat și acesta este principalul motiv pentru care RayTracing nu este folosit la scară largă în simulările 3D computerizate ce trebuie să ruleze în timp real. Deși RayTracing-ul este intensiv computațional, existența independenței între calculele efectuate în cadrul algoritmului îl fac un candidat optim pentru paralelizare.

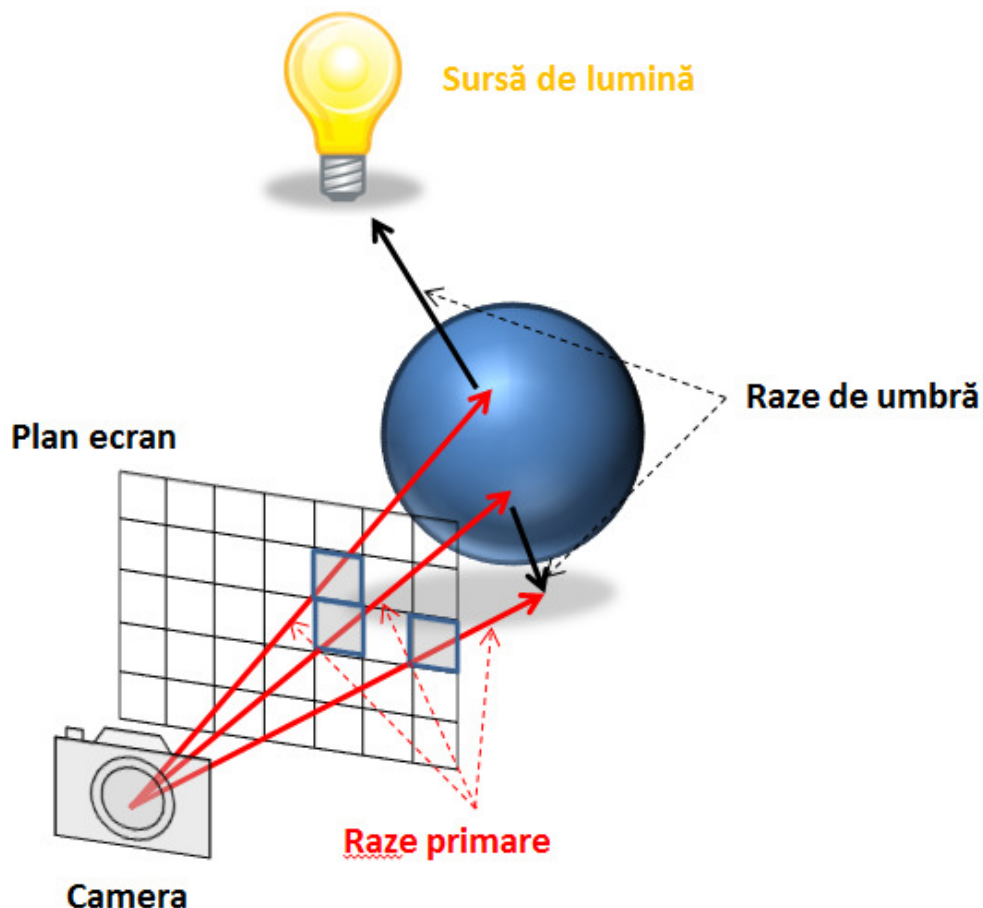
#### **4.2.1.1 Descrierea algoritmului RayTracing**

RayTracing folosește pentru obținerea imaginii 2D finale care va fi afișată pe ecran, un set de raze care pornesc din punctul din care se privește scena. Razele vor trece prin fiecare pixel al planului virtual al ecranului urmând să intersecteze obiectele din scena 3D. Aceste raze se numesc raze primare iar procedeul prin care ele sunt folosite se numește raycasting reprezentând primă etapă a algoritmului RayTracing.

Imaginea 2D din planul virtual al ecranului corespunde unei ferestre de proiecție, iar pentru determinarea direcțiilor razelor, este necesară calcularea pozițiilor pixelilor în spațiul scenei, lucru care se realizează printr-o transformare de coordonate.

După ce se calculează direcțiile pentru fiecare rază, urmează etapa a doua a algoritmului : intersecția acestora cu primitivele din scenă. Astfel, pentru fiecare rază primară, se determină primitiva cea mai apropiată (care este intersectată prima), calculându-se poziția punctului de intersecție.

După ce s-a obținut punctul de intersecție, se determină vizibilitatea față de sursele de lumină, pentru a se afla dacă punctul este sau nu umbrat. Pentru aceasta, se generează pentru fiecare sursă de lumină câte o rază ce are ca origine punctul de intersecție și are direcția către sursa de iluminare. Aceste raze se numesc raze de umbră și sunt intersectate cu primitivele din scenă pentru a determina dacă există obstacole între punctul de intersecție și sursele de lumină. La prima intersecție găsită, căutarea este întreruptă deoarece lumina de la sursă va fi blocată și nu va mai ajunge la punctul de pe obiect și astfel sursa de lumină către care este îndreptată raza de umbră curentă nu mai este luată în considerare la calcularea culorii punctului.



*Figura 4-10: Algoritm RayTracing*

Culorile pixelilor corespunzători punctelor de intersecție se determină folosind calcule ce implementează un model de iluminare, la calculul culorii finale contribuind și rezultatele obținute din razele de reflexie și refracție. După ce se efectuează toate aceste operații pentru fiecare pixel al imaginii din planul virtual al ecranului, fiecare pixel va avea atribuită o culoare iar imaginea generată poate fi afișată.

Datorită faptului că operațiile executate pentru fiecare pixel pot fi executate independent unele față de celelalte, există un potențial ridicat de paralelism pentru algoritm fiind posibilă executarea concurrentă a calculelor pentru determinarea culorii finale a fiecărui pixel al imaginii.

Pseudo-codul pentru algoritm este următorul [MOL96] :

```
Pentru fiecare pixel al imaginii
{
    Construiește rază de la observator prin pixel
    Inițializează CelMaiApropiatT cu infinit și
    CelMaiApropiatObiect cu NULL

    Pentru fiecare obiect al scenei
    {
        Dacă Există Intersecție(rază, obiect)
        {
            Dacă t intersecție < CelMaiApropiatT
            {
                CelMaiApropiatT = intersecție
                CelMaiApropiatObiect = obiect curent
            }
        }
    }
    Dacă CelMaiApropiatObiect este NULL
    {
        culoare pixel = culoarea de fundal
    }
    Altfel
    {
        Construiește raze de umbră către toate sursele
        de lumină din scenă

        Dacă intersecția nu este în umbră
        {
            Dacă suprafața generează reflexie
            {
                Construiește raza de reflexie
                Reapelează algoritm pentru noua rază
            }

            Dacă suprafața este transparentă
            {
                Construiește raza de refracție
                Reapelează algoritm pentru noua rază
            }
            culoare pixel = culoare calculată folosind
            model de iluminare
        }
    }
    Altfel
    {
```

```
        culoare pixel = culoare calculată folosind
        model de iluminare și faptul ca pixelul
        curent este în umbră
    }
}
}
```

Se observă în pseudocodul de mai sus, faptul ca algoritmul este unul recursiv (prin reapelarea algoritmului pentru razele de reflexie și refracție), pentru fiecare rază primară fiind generate alte raze, denumite raze secundare, care vor contribui la culoarea finală a pixelului. Astfel, în practică se folosește o adâncime maximă până la care se apelează algoritmul (numărul maxim de reflexii/refracții secundare).

Razele de reflexie și refracție nu sunt generate decât dacă obiectul intersectat are un material care permite reflexia respectiv refracția (grad de transparență). Etapa cea mai importantă și costisitoare este cea a testelor de intersecție deoarece conține calcule cu vectori.

Razele generate au următoarea formă parametrizată:

$$\vec{r}(t) = \vec{o} + t\vec{d}$$

unde :

- $\vec{o}$  : originea razei
- $\vec{d}$  : direcția razei
- $t$  : distanța parcursă de raza.

În pseudocodul prezentat,  $t$  reprezintă acest parametru și are valoarea distanței pe care o rază trebuie să o parcurgă până la punctul de intersecție.

#### ***4.2.1.2 Calculul culorii pixelului folosind un model de iluminare***

Exista mai multe modele de iluminare (Phong [PHO75], Gouraud, Blinn [BLI77], Lambert, etc) care pot fi folosite pentru a determina modul în care lumina afectează obiectele din scenă.

Pentru a obține astfel culoarea unui obiect iluminat vom avea următoarele componente :

- Componenta emisivă

- Componenta ambientală
- Componenta difuză
- Componenta speculară

Contribuția fiecărei dintre aceste componente este calculată ca o combinație dintre proprietățile de material ale obiectului (factorul de strălucire, culoarea materialului) și proprietățile sursei de lumină (culoarea sursei de lumină, poziția sursei de lumină).

Astfel pentru culoarea finală asociată obiectului vom avea :

$$\text{culoareObiect} = \text{emisiva} + \text{ambientala} + \text{difuza} + \text{speculara}$$

### **Componenta emisivă**

Aceasta reprezintă lumina emisă de un obiect și nu ține cont de nici o sursă de lumină. Dacă un obiect care are o anumită culoare emisivă s-ar afla într-o scenă complet întunecată atunci el ar apărea exact cu această culoare.

O utilizare des întâlnită pentru componenta emisivă este aceea de a simula strălucirea unui obiect.

Avem astfel :

$$\text{emisiva} = Ke$$

$Ke$  – culoarea emisivă a materialului

### **Componenta ambientală**

Aceasta reprezintă lumina care a fost reflectată de către obiectele din scena de atât de multe ori încât pare să vină de peste tot.

Astfel, lumina ambientală nu vine dintr-o direcție anume, apărând ca venind de fapt din toate direcțiile. Din aceasta cauză, componenta ambientală este independentă de poziția sursei de lumină.

Componenta ambientală depinde de culoarea de material ambientală a obiectului și de culoarea ambientală a luminii.

Similar componentei emisive, componenta ambientală este o constantă (se poate extinde modelul atribuind fiecărei lumini din scena o culoare ambientală proprie).

Avem astfel :

$$\text{ambientala} = Ka * \text{culoareaAmbientalaGlobala}$$

$Ka$  – culoarea de material ambientală a obiectului

*culoareaAmbientalaGlobala* – culoarea ambientală a luminii

### Componenta difuză

Aceasta reprezintă lumina reflectată de suprafața obiectului în mod egal în toate direcțiile.

Cantitatea de lumina reflectată este proporțională cu unghiul de incidență al razei de lumina cu suprafața obiectului.

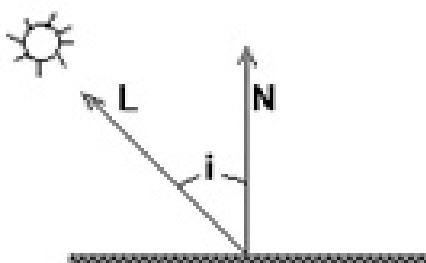


Figura 4-11: Calcul componentă difuză

Avem astfel :

$$difuză = Kd * culoareLumina * \max(N \cdot L, 0)$$

*Kd* - culoarea de material difuză a obiectului

*culoareLumina* – culoarea difuză a luminii

*N* – normala la suprafață (normalizată)

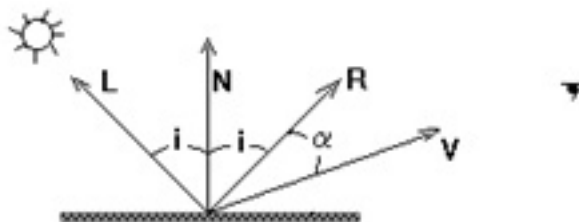
*L* – vectorul direcției luminii incidente (normalizat)

$\max(N \cdot L, 0)$  – produsul scalar  $N \cdot L$  ce reprezintă măsura unghiului dintre acești 2 vectori; astfel, dacă  $i$  este mai mare decât  $\pi/2$  valoarea produsului scalar va fi mai mică decât 0 acest lucru însemnând că suprafața nu primește lumina (sursa de lumină se află în spatele suprafeței) și de aici și formula care asigură faptul că în acest caz suprafața nu primește lumina difuză

### Componenta speculară

Un reflector perfect, de exemplu o oglindă, reflectă lumina numai într-o singură direcție *R*, care este simetrică cu *L* față de normala la suprafață. Deci numai un observator situat exact pe direcția respectivă va percepe raza reflectată.



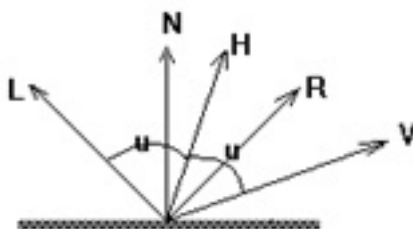


*Figura 4-12: Calcul componentă speculară*

Componenta speculară reprezintă lumina reflectată de suprafața obiectului numai în jurul acestei direcții R.

În modelul Phong se aproximează scăderea rapidă a intensității luminii reflectate atunci când  $\alpha$  crește prin  $\cos(\alpha)^n$ , unde  $n$  este exponentul de reflexie speculară al materialului (shininess).

O altă formulare a modelului Phong se bazează pe vectorul median, notat cu H. El face unghiuri egale cu L și cu V. Dacă suprafața ar fi orientată astfel încât normala să aibă direcția lui H, atunci observatorul ar percepe lumina speculară maximă (deoarece ar fi pe direcția razei reflectate specular).



*Figura 4-13: Calcul componentă speculară folosind vectorul median*

Termenul care exprimă reflexia speculară este în acest caz :

$$(N \cdot H)^n \text{ unde } H = (L + V) \text{ (normalizat)}$$

Atunci când sursa de lumină și observatorul sunt la infinit, utilizarea termenului  $N \cdot H$  este avantajoasă deoarece  $H$  este constant.

După cum se observă, față de celelalte 3 componente, componenta speculară depinde și de poziția observatorului. Dacă observatorul nu se află într-o poziție unde poate vedea razele reflectate atunci nu va vedea reflexie speculară pentru zona respectivă. De asemenea nu va vedea reflexie speculară dacă lumina se află în spatele suprafeței.

Ținând cont de toate acestea avem pentru componenta speculară următoarea formulă :

$$speculara = K_s * culoareLumina * primesteLumina * (max(N \cdot H, 0))^n$$

$K_s$  - culoarea de material speculară a obiectului

$culoareLumina$  – culoarea speculară a luminii

$N$  – normala la suprafață (normalizată)

$L$  – vectorul direcției luminii incidente (normalizat)

$H$  – vectorul median (normalizat)

$primesteLumina$  – 1 dacă  $N \cdot L$  este mai mare decât 0, 0 în caz contrar

### Atenuarea intensității luminii

Atunci când sursa de lumină punctiformă este suficient de îndepărtată de obiectele scenei vizualizate, vectorul  $L$  este același în orice punct. Sursa de lumină este numită în acest caz direcțională. Aplicând modelul pentru vizualizarea a două suprafețe paralele construite din același material, se va obține o aceeași intensitate (unghiul dintre  $L$  și normală este același pentru cele două suprafețe). Dacă proiecțiile suprafețelor se suprapun în imagine, atunci ele nu se vor distinge. Aceasta deoarece în model nu se ține cont de faptul că intensitatea luminii descrește proporțional cu inversul pătratului distanței de la sursa de lumină la obiect. Deci, obiectele mai îndepărtate de sursă sunt mai slab luminate. O posibilă corecție a modelului, care poate fi aplicată pentru surse poziționale (la distanță finită de scena) este:

$$culoareObiect = emisiva + ambientala + factorAtenuare * (difuza + speculara)$$

$factorAtenuare = 1/d^2$  este o funcție de atenuare;  $d$  este distanța de la sursă la punctul de pe suprafață considerat.

Corecția nu satisface cazurile în care sursa este foarte îndepărtată; de asemenea, dacă sursa este la distanță foarte mică de scenă, intensitățile obținute pentru două suprafețe cu același unghi  $i$ , între  $L$  și  $N$ , vor fi mult diferite.

O aproximare mai bună este următoarea:

$$factorAtenuare = 1/(K_c + K_l * d + K_q * d^2)$$

unde  $K_c$  este factorul de atenuare constant,  $K_l$  este factorul de atenuare liniar și  $K_q$  este factorul de atenuare pătratic

## Lumina spot

Un efect interesant este acela de a face lumina să aibă un efect de spot, astfel încât lumina să afecteze doar obiectele ce se află în interiorul unui con de lumină (similar cu efectul produs de exemplu de o lampă de birou).

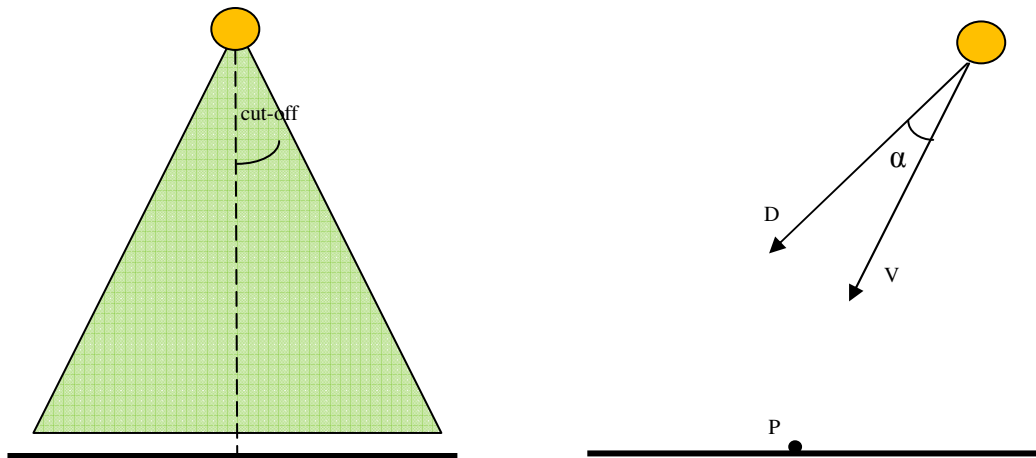


Figura 4-14: *Lumină de tip spot*

Pentru a crea o astfel de lumină este nevoie de următorii parametrii :

- poziția luminii
- direcția spotului (D)
- unghiul de cut-off al spotului (controlează “deschiderea” conului de lumină după cum se vede din figura de mai sus)
- poziția punctului care se dorește a fi iluminat (P)

Astfel din aceste informații se poate deduce și V care reprezintă vectorul direcției de la lumina la punctul P.

Pentru a determina dacă punctul P se află sau nu în conul de lumină este de ajuns să determinăm dacă este îndeplinită condiția :

$$\cos(\alpha) \geq \cos(\text{cut-off}) , \text{ unde } \cos(\alpha) = V \cdot D$$

### 4.2.1.3 Testele de intersecție

Obiectele din scena 3D care sunt folosite în implementarea algoritmului sunt descrise analitic. În continuare sunt prezentate ecuațiile prin care sunt descrise aceste obiecte precum și rezultatele intersecției lor cu razele de lumină.

#### 4.2.1.3.1 Sfera

Fie o sferă de centru  $\mathbf{c} = (c_x, c_y, c_z)$  și rază  $r$  ce este descrisă prin ecuația implicită:

$$(x - c_x)^2 + (y - c_y)^2 + (z - c_z)^2 - r^2 = 0$$

Această ecuație se poate scrie simplificat în formă vectorială:

$$(\mathbf{p} - \mathbf{c}) \cdot (\mathbf{p} - \mathbf{c}) - r^2 = 0$$

În această ecuație  $\mathbf{p} = (x, y, z)$  este un punct de pe sferă. Pentru a afla intersecția cu raza, se consideră  $\mathbf{p}$  aparținând razei descrisă de ecuația parametrică  $\mathbf{p}(t) = \mathbf{o} + t\mathbf{d}$ .

Se obține astfel:

$$(\mathbf{o} + t\mathbf{d} - \mathbf{c}) \cdot (\mathbf{o} + t\mathbf{d} - \mathbf{c}) - r^2 = 0$$

Dezvoltând, rezultă ecuația pătratică:

$$(\mathbf{d} \cdot \mathbf{d})t^2 + [2(\mathbf{o} - \mathbf{c}) \cdot \mathbf{d}]t + (\mathbf{o} - \mathbf{c}) \cdot (\mathbf{o} - \mathbf{c}) - r^2 = 0$$

cu delta  $\Delta = 2\sqrt{((\mathbf{o} - \mathbf{c}) \cdot \mathbf{d})^2 - (\mathbf{d} \cdot \mathbf{d})(\mathbf{o} - \mathbf{c}) \cdot (\mathbf{o} - \mathbf{c}) - r^2}$ )}

Din valoarea lui delta se poate afla câte intersecții sunt între rază și sferă :

- delta > 0, există 2 intersecții
- delta = 0 există o singură intersecție
- delta < 0 raza nu se intersectează cu sfera.

Un caz special este atunci când raza intersectează sfera în două locuri, punctul de intersecție corespunzând celui mai mic  $t$  pozitiv.

$$t = \frac{-(\mathbf{o} - \mathbf{c}) \cdot \mathbf{d} \pm \sqrt{\Delta}}{\mathbf{d} \cdot \mathbf{d}}$$

Normala la sferă în punctul de intersecție este dată de gradientul ecuației vectoriale

$$\mathbf{N} = 2(\mathbf{p} - \mathbf{c})$$

#### 4.2.1.3.2 Planul

Fie un plan versorul normalei  $\mathbf{n} = (A, B, C)$  aflat la distanța  $D$  față de origine. Planul poate fi descris astfel prin ecuația implicită

$$Ax + By + Cz + D = 0$$

Ecuția se poate scrie simplificat în formă vectorială:

$$\mathbf{n} \cdot \mathbf{p} + D = 0$$

unde  $\mathbf{p} = (x, y, z)$  este un punct din plan.

Pentru a determina intersecția cu raza, se consideră  $\mathbf{p}$  aparținând razei descrisă de ecuația parametrică  $\mathbf{p}(t) = \mathbf{o} + t\mathbf{d}$ .

Înlocuind, se obține:

$$\mathbf{n} \cdot (\mathbf{o} + t\mathbf{d}) + D = 0$$

Se observă că este o ecuație liniară din care se poate scoate parametrul  $t$

$$t = -\frac{D + \mathbf{n} \cdot \mathbf{o}}{\mathbf{n} \cdot \mathbf{d}}$$

Dacă numitorul este 0, atunci raza este paralelă la plan, iar dacă atât numitorul cât și numărătorul sunt 0 atunci raza este inclusă în plan, ecuația fiind satisfăcută pentru orice  $t$ .

De interes este doar cazul în care numitorul este pozitiv, punctul de intersecție fiind dat de parametrul  $t$  pozitiv. Normala la plan în punctul de intersecție este dată de gradientul ecuației implicite și valoarea este în acest caz chiar  $\mathbf{n}$ .

#### 4.2.1.3.3 Triunghiul

Un triunghi este definit de trei puncte  $\mathbf{V}_1$ ,  $\mathbf{V}_2$  și  $\mathbf{V}_3$ , care dacă nu sunt coliniare, definesc un plan. Cel mai facil mod pentru descrierea acestui plan este folosind coordonatele baricentrice.

$$\mathbf{p}(u, v) = (1 - u - v)\mathbf{V}_0 + u\mathbf{V}_1 + v\mathbf{V}_2$$

Dacă se îndeplinesc condițiile  $0 \leq u \leq 1$  și  $0 \leq v \leq 1$  atunci punctul  $\mathbf{p}(u, v)$  aparține triunghiului.

Pentru a afla intersecția cu raza, se consideră  $\mathbf{p}$  aparținând razei descrisă de ecuația parametrică  $\mathbf{p}(t) = \mathbf{o} + t\mathbf{d}$ . Înlocuind, se obține:

$$\mathbf{o} + t\mathbf{d} = (1 - u - v)\mathbf{V}_0 + u\mathbf{V}_1 + v\mathbf{V}_2$$

Rearanjând termenii se obține sistemul liniar de ecuații:

$$\begin{bmatrix} -\mathbf{d} & \mathbf{V}_1 - \mathbf{V}_0 & \mathbf{V}_2 - \mathbf{V}_0 \end{bmatrix} \begin{bmatrix} t \\ u \\ v \end{bmatrix} = \mathbf{o} - \mathbf{V}_0$$

Rezolvând acest sistem se obțin valorile pentru  $t$ ,  $u$  și  $v$ . Se fac notațiile  $\mathbf{E}_1 = \mathbf{V}_1 - \mathbf{V}_0$ ,  $\mathbf{E}_2 = \mathbf{V}_2 - \mathbf{V}_0$  și  $\mathbf{T} = \mathbf{o} - \mathbf{V}_0$ . Soluția ecuației anterioare se poate găsi cu regula lui Cramer:

$$\begin{bmatrix} t \\ u \\ v \end{bmatrix} = \frac{1}{\begin{vmatrix} -\mathbf{d}, \mathbf{E}_1, \mathbf{E}_2 \end{vmatrix}} \begin{bmatrix} \begin{vmatrix} -\mathbf{T}, \mathbf{E}_1, \mathbf{E}_2 \end{vmatrix} \\ \begin{vmatrix} -\mathbf{d}, \mathbf{T}, \mathbf{E}_2 \end{vmatrix} \\ \begin{vmatrix} -\mathbf{d}, \mathbf{E}_1, \mathbf{T} \end{vmatrix} \end{bmatrix} = \frac{1}{(\mathbf{d} \times \mathbf{E}_2) \cdot \mathbf{E}_1} \begin{bmatrix} (\mathbf{T} \times \mathbf{E}_1) \cdot \mathbf{E}_2 \\ (\mathbf{d} \times \mathbf{E}_2) \cdot \mathbf{T} \\ (\mathbf{T} \times \mathbf{E}_1) \cdot \mathbf{d} \end{bmatrix}$$

Se consideră intersecție doar dacă  $t$  obținut este pozitiv. Normala în punctul de intersecție este dată de interpolarea normalelor la vârfuri  $\mathbf{N}_1$ ,  $\mathbf{N}_2$ ,  $\mathbf{N}_3$  în funcție de coordonatele  $(u, v)$  ale punctului de intersecție.

$$\mathbf{N} = (1 - u - v)\mathbf{N}_1 + u\mathbf{N}_2 + v\mathbf{N}_3$$

#### 4.2.2 RAYTRACING FOLOSIND TEHNICI GPGPU

În ultimii ani, interesul pentru a oferi soluții de redare ce folosesc RayTracing a înregistrat o creștere semnificativă atât în industria producătorilor de hardware (GPU și CPU) cât și în cea a dezvoltatorilor de jocuri. Există astfel în prezent modificări ale motoarelor grafice din unele jocuri apărute pentru a folosi RayTracing ([QUA04], [QUA06], [QWR08], [WOL10]), programe demonstrative de la producătorii de GPU ([NVI09a]), încercări de a implementa motoare grafice ce folosesc RayTracing de la dezvoltatorii de jocuri ([CAR11]), etc.

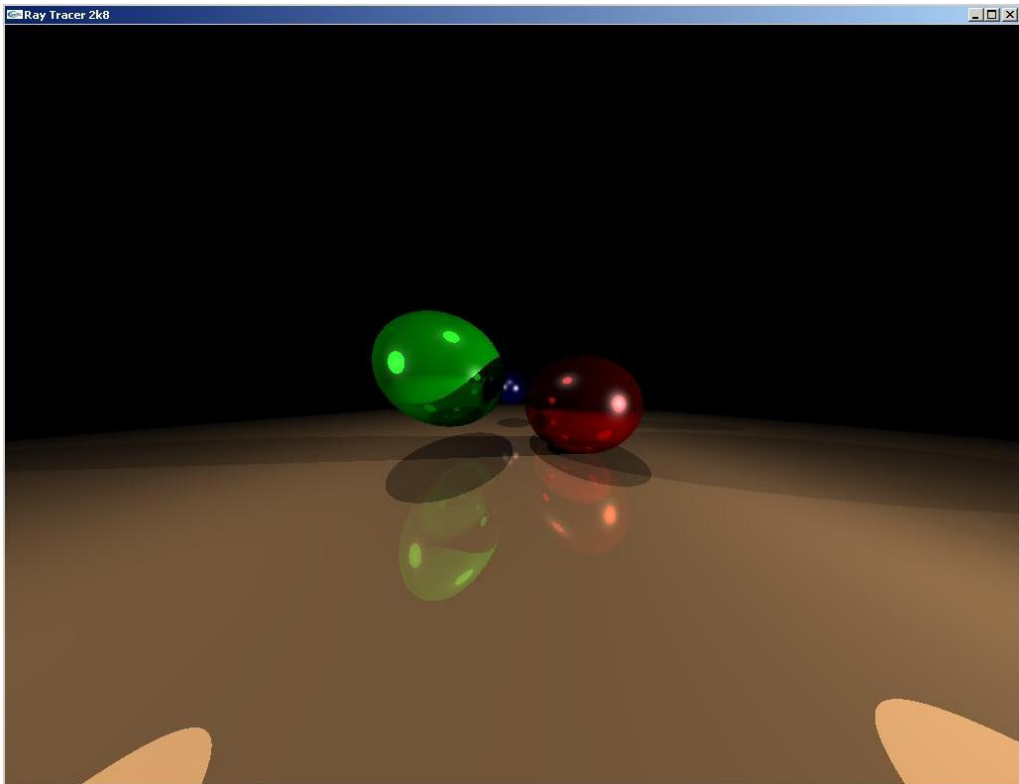
Soluțiile/implementările curente de RayTracing pe GPU se concentrează în principal pe următoarele aspecte :

- Optimizarea algoritmului RayTracing folosind diverse tehnici de partiționare spațială pentru reprezentarea scenei 3D folosind arbori K-D ([FOL05]), BVH (Volume Încadratoare Ierarhice), etc;
- Execuția algoritmului RayTracing pe GPU folosind banda grafică; acest lucru presupune de obicei scrierea de *shadere* (vertex shader, geometry shader și fragment shader);

- Utilizarea de abordări hibride, unde părți ale algoritmului sunt executate pe GPU folosind shadere;
- Abordări distribuite ce executa algoritmul RayTracing pe mai multe mașini folosind un framework specializat ([WAL04]).

Spre deosebire de soluțiile prezentate mai sus, soluția descrisă în continuare își propune să execute algoritmul RayTracing în totalitate pe procesoarele scalare ale GPU utilizând o implementare GPGPU, fiind astfel o abordare ce se diferențiază de utilizarea standard a GPU, care folosește banda grafică cu programe de tip shader.

Având în vedere că soluția propusă se concentrează pe implementarea algoritmului direct pe procesoarele scalare ale GPU am decis sa fie folosită o varianta a algoritmului RayTracing care este intensiv computațională, fără alte optimizări suplimentare. În acest fel, se pot compara rezultatele obținute de variantele GPGPU (single și multi GPGPU) și de implementarea doar pe CPU.



***Figura 4-15: Scena redată folosind RayTracing***

### 4.2.3 ADAPTARE ALGORITM RAYTRACING PENTRU GPGPU

Algoritmul pentru împărțirea calculului RayTracing pe firele de execuție CUDA este următorul :

- Datele inițiale ale scenei sunt făcute disponibile în memoria globală a GPU;
- Setul de raze este determinat de o matrice de pixeli ce are dimensiunea ecranului (Figura 3-16);
- Matricea ecran va fi la rândul ei împărțită în sub-matrice ce au dimensiuni fixe (8,16 sau 32);
- Fiecare sub-matrice va fi procesată de un bloc de fire de execuție unde fiecare fir de execuție va fi responsabil de calculele executate pentru un element al sub-matricei (un pixel și raza care trebuie să treacă prin el);
- Pentru implementarea multi-GPGPU, imaginea se va împărți la numărul de GPU, fiecare GPU având asignată o porțiune din imagine (de exemplu, dacă se folosesc două GPU, GPU0 va procesa jumătatea superioară a ecranului și GPU1 va procesa jumătatea inferioară a ecranului)

```
GPUWorker gpu0(0);  
GPUWorker gpu1(1);  
  
gpu0.callAsync(bind(launch_thread,pixels,  
device,dimGrid,dimBlock));  
  
gpu1.callAsync(bind(launch_thread,pixels,  
device,dimGrid,dimBlock));
```

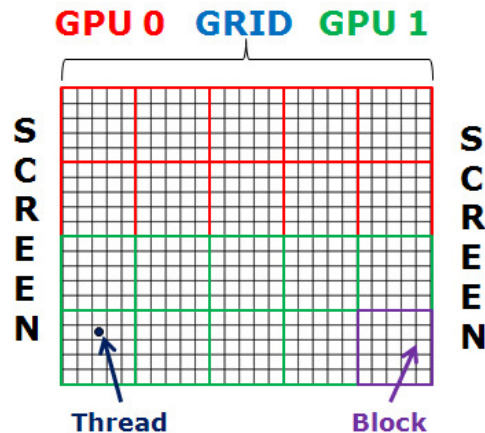


Figura 4-16: Împărțirea pe GPU / blocuri SIMD a sarcinilor

Pentru implementarea *single-GPGPU* pseudocodul algoritmului este următorul:  
:



```
DimensiuneBlocFireExecuție = (32,32)
DimensiuneGridFireExecuție =
(LățimeRezoluțieScenă / dimensiuneXbloc,
ÎnălțimeRezoluțieScenă / dimensiuneYbloc)

Copiază date în memoria globală GPU de pe mașina gazdă

Lansează în paralel pe GPU funcție kernel(dimensiune
bloc,dimensiune grid)

Sincronizează prin barieră terminarea firelor de execuție

Copiază date din memoria globală GPU pe mașina gazdă
```

Funcția kernel ce se execută pe fiecare fir de execuție de pe GPU este :

```
X = IdBlockX * DimensiuneBlocX + IdFirDeExecuțieX
Y = IdBlockY * DimensiuneBlocY + IdFirDeExecuțieY

IndiceMemorieImagine = X * LățimeRezoluțieScenă + Y

RayTracing(pixel ecran identificat de IndiceMemorieImagine)
```

Pentru implementarea *multi-GPGPU* pseudocodul algoritmului se modifică astfel :

```
DimensiuneBlocFireExecuție = (32,32)
DimensiuneGridFireExecuție =
(LățimeRezoluțieScenă/dimensiuneXbloc/nrGPU,
ÎnălțimeRezoluțieScenă/dimensiuneYbloc/nrGPU)

Pentru fiecare GPU
{
    Copiază date în memoria globală GPU de pe mașina gazdă

    Lansează în paralel pe GPU funcție kernel(dimensiune
    bloc,dimensiune grid, idGPU)

    Sincronizează prin barieră terminarea firelor de execuție

    Copiază date din memoria globală GPU pe mașina gazdă
```

}

Funcția kernel ce se execută pe fiecare fir de execuție de pe GPU este :

```
X = IdBlockX * DimensiuneBlockX + IdFirDeExecuțieX
Y = IdBlockY * DimensiuneBlockY + IdFirDeExecuțieY

Offset = ((LățimeRezoluțieScenă*ÎnălțimeRezoluțieScenă)/nrGPU)
         * idGPU
IndiceMemorieImagine = X * LățimeRezoluțieScenă + Y + Offset
RayTracing(pixel ecran identificat de Indice)
```

#### 4.2.4 ANALIZĂ TIMP DE EXECUȚIE

Timpul de execuție al algoritmului RayTracing în implementarea pe CPU este următorul :

$$T( W*H*NrObiecte*(NivelRazeSecundare+1)),$$

unde,

W – lățimea în pixeli a rezoluției scenei

H – înălțimea în pixeli a rezoluției scenei

NrObiecte – numărul de obiecte din scenă

NivelRazeSecundare – numărul maxim admis de ricoșee ce pot genera raze secundare în urma unor reflexii sau refracții

Timpul de execuție al algoritmului RayTracing în implementarea *single-GPGPU* este următorul:

$$T((( W*H / (nrSM * nrSP) ) * (NrObiecte * (NivelRazeSecundare + 1))))$$

+ *constantaOverhead*,

unde,

W – lățimea în pixeli a rezoluției scenei

H – înălțimea în pixeli a rezoluției scenei

NrObiecte – numărul de obiecte din scenă

NivelRazeSecundare – numărul de ricoșee ce pot genera raze secundare în urma unor reflexii sau refracții

nrSM – numărul de multiprocesoare de pe GPU

nrSP – numărul de procesoare scalare de pe fiecare multiprocesor  
constantaOverhead – constantă introdusă de operațiile de transfer de memorie și sincronizare

Timpul de execuție al algoritmului RayTracing în implementarea *multi-GPGPU* este următorul:

$$T( ((W*H / (nrSM*nrSP*nrGPU)) * (NrObiecte * (NivelRazeSecundare + 1))) + nrGPU * constantaOverhead,$$

unde,

W – lățimea în pixeli a rezoluției scenei

H – înălțimea în pixeli a rezoluției scenei

NrObiecte – numărul de obiecte din scenă

NivelRazeSecundare – numărul de ricoșee ce pot genera raze secundare în urma unor reflexii sau refracții

nrSM – numărul de multiprocesoare de pe GPU

nrSP – numărul de procesoare scalare de pe fiecare multiprocesor

nrGPU – numărul de GPU-uri disponibile

constantaOverhead – constantă introdusă de operațiile de transfer de memorie și sincronizare

#### 4.2.5 EVALUAREA REZULTATELOR

Pentru scena de test au fost folosite :

- O rezoluție de 1024x768
- Iluminare difuză și speculară de la 3 surse de lumină
- Generarea de umbre (care este obținută natural de RayTracing)
- Generarea de reflexii, cu o limită de maxim 3 raze secundare ce pot rezulta în urma reflexiilor.

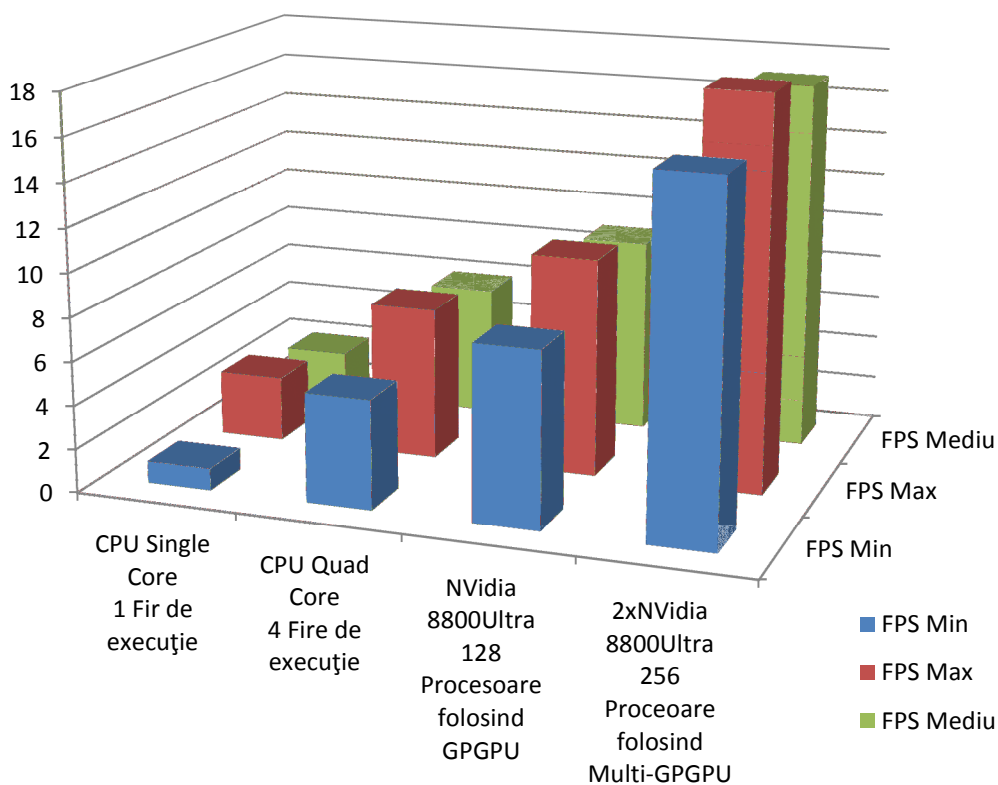
Pentru implementarea pe CPU (Single Core – 1 fir de execuție și Quad Core – 4 fire de execuție) a fost utilizat un CPU Intel Core2Quad Q6600 și algoritmul RayTracing a fost rulat folosind 1 fir de execuție, respectiv 4 fire de execuție.

Pentru implementările GPGPU (single și multi) au fost folosite 2 GPU-uri Nvidia 8800 Ultra montate în același calculator și toolkit-ul CUDA 2.2.

Au fost testate toate cele 3 implementări și rezultatele obținute sunt cele prezentate în continuare.

	FPS Min	FPS Max	FPS Mediu
CPU Single Core 1 fir de execuție	1	3	2.083
CPU Quad Core 4 fire de execuție	5	7	5.983
NVidia 8800Ultra 128 procesoare folosind GPGPU	8	10	8.967
2xNVidia 8800Ultra 256 procesoare folosind GPGPU	16	18	16.980

***Tabel 4-1: Rezultate implementare RayTracing folosind GPGPU***



***Figura 4-17: Rezultate rulare algoritm RayTracing***

Rezultatele obținute demonstrează faptul că algoritmul RayTracing scalează foarte bine atunci când este executat paralel și ca se obține un spor de performanță prin implementarea GPGPU și Multi-GPGPU.

## **4.3 UTILIZAREA GPGPU ÎN CADRUL SERVERELOR MMO**

### **4.3.1 IDENTIFICARE OPERAȚII CE POT FI PARALELIZATE**

#### **4.3.1.1 Descriere generală**

Cantitatea de informații pe care o aplicație trebuie să le proceseze este în continuă creștere. Aplicațiile MMO monitorizează și procesează un număr de peste 1 milion ([LAK09]) de operații pe secundă în perioadele de maximă încărcare.

Este clar că scalabilitatea în ceea ce privește numărul de evenimente procesate este o cerință importantă a aplicațiilor MMO.

Primele generații de sisteme de procesare a evenimentelor în cadrul aplicațiilor MMO au fost de tip centralizat, însă de curând au apărut soluții ce încearcă să obțină scalabilitatea folosind o abordare distribuită.

Partiționarea proceselor se dovedește însă a fi destul de dificilă datorită numeroaselor dependențe ce pot apărea între procese cât și între resursele ce le procesează. De obicei, partiționarea se face într-un mod vertical, existând un singur set de partiționare ce folosește un mecanism de distribuire centralizat.

#### **4.3.1.2 Logica spațiului virtual**

Nucleul unei aplicații MMO este reprezentat de o buclă internă în cadrul simulării, care se execută sincronizată cu un interval de timp ce corespunde unui anumit număr de frame-rate (de obicei 30 sau 60 de FPS - cadre pe secundă).

Din punct de vedere conceptual, starea lumii virtuale a unei aplicații MMO poate fi văzută sub forma unei tabeli de date ce conține obiectele spațiului virtual incluzând atât utilizatorii precum și obiectele cu care aceștia interacționează la un anumit moment de timp.

Pentru a face față cerințelor stricte de simulare în timp real, starea activă a spațiului virtual este de obicei păstrată în memoria serverelor centrale ce asigură funcționarea acestuia. Discurile fizice sunt folosite de obicei pentru a asigura persistența sau pentru a stoca informații auxiliare.

Pe parcursul fiecărei iterații ( *tick* ) a buclei principale a simulării, porțiuni ale stării sunt actualizate în concordanță cu logica simulării spațiului virtual. Aceste actualizări pot fi declanșate de acțiuni ale utilizatorilor, expirarea unui anumit interval de timp sau alte evenimente.

Este posibil ca o acțiune singulară a unui utilizator să declanșeze actualizări fizice multiple ale stării curente. De exemplu, o comandă de deplasare la nivelul utilizatorului se poate transla într-o serie de mici actualizări ale poziției acestuia care se desfășoară pe parcursul a mai multor *tick-uri*.

Vor fi enumerate în continuare câteva dintre principalele operații ce se execută la nivelul logicii simulării spațiului virtual și care se execută de obicei folosind un mecanism distribuit, fiind astfel susceptibile la posibile optimizări.

#### **4.3.1.3 Fizica spațiului virtual**

Pentru a asigura consistența și corectitudinea funcționării spațiului virtual, aplicațiile MMO trebuie să execute, la nivelul serverului, un număr foarte mare de operații pentru a simula o fizică realistă a spațiului virtual (detectia coliziunilor, forțe de frecare, forța de gravitație, etc).

Acestea reprezintă gama de operații care sunt cele mai frecvent efectuate de către serverele spațiului virtual și astfel sunt principalele consumatoare de timp de calcul.

Spațiul virtual trebuie să asigure corectitudinea interacțiunilor fizice între entitățile care îl populează și astfel trebuie să execute foarte multe calcule de detecție a coliziunilor. De obicei, la fiecare deplasare / actualizare a poziției unui caracter în spațiul virtual trebuie asigurat faptul că acesta are voie să ajungă la poziția respectivă și că nu intră, de exemplu, în zidul unei clădiri.

Deoarece calculele de detecție a coliziunilor sunt foarte costisitoare, în practică se limitează numărul entităților între care se efectuează aceste calcule (caracterele cu terenul, caracterele cu anumite cladiri, etc) și sunt folosite de obicei modele de aproximare ale obiectelor asupra cărora se execută calculele.

#### **4.3.2 IMPLEMENTAREA SIMULARILOR FIZICII FOLOSIND GPGPU**

Așa cum s-a precizat și anterior, operațiile cele mai costisitoare din punct de vedere al timpului de calcul care sunt executate de către serverele spațiilor virtuale 3D MMO sunt cele de detecție a coliziunilor. În practică, se utilizează pentru detecția coliziunilor metode ce iau în considerație :

- fie poligoanele din care sunt alcătuite obiectele 3D
- fie o aproximare a obiectelor 3D folosind volume încadratoare

În mod uzual se folosește cea de-a doua metodă, deoarece prin simplificarea suprafețelor pentru care trebuie facute calculele de coliziune se obține un timp mai mic de execuție.

În continuare se propune o modalitate de implementare a acestor tipuri de operații folosind tehnici GPGPU.

#### **4.3.2.1 Volum încadratoare**

În grafica pe calculator, un volum încadrator pentru un set de obiecte este un volum închis care conține toate obiectele mulțimii respective. Acesta este folosit pentru a îmbunătăți eficiența operațiilor geometrice asupra unor forme complexe, folosind forme geometrice mai simple. Astfel de volume mai simple permit folosirea unor metode mai eficiente de testare a coliziunilor.

Volumele încadratoare cele mai folosite sunt:

- Sfere
- Elipsoizi
- Cilindri
- Paralelipipede

Sunt folosite următoarele metode de detecție a coliziunilor folosind volume încadratoare:

##### *A) AABB*

AABB (Axis Aligned Bounding Box) este un paralelipiped încadrator ale cărui fețe sunt paralele cu axele sistemului de coordonate.

Pentru a determina dacă două AABB se intersectează, se testează separat pentru fiecare dintre axele  $x$ ,  $y$ ,  $z$ . Acest test separat pentru fiecare din cele 3 axe va determina coliziunea a două AABB chiar și în cazul în care unul dintre ele îl conține pe celălalt.

Cele două AABB a căror intersecție este verificată, sunt specificate prin coordonatele centrului  $C$  și un vector  $E$ , care reprezintă distanța de la centru la

fiecare față a paralelipipedului, rezultând 3 valori, corespunzătoare fețelor paralele cu axele  $x, y, z$ .

Se calculează pentru fiecare axă distanța coordonatelor centrelor celor două AABB, obținându-se vectorul T.

$$T = AABB1.C - AABB2.C$$

Se compară distanțele corespunzătoare celor 3 axe cu suma distanțelor D ale celor 2 AABB. Astfel cele două AABB se intersectează atunci când:

$$|T.x| \leq AABB1.E.x + AABB2.E.x \ \&\&$$

$$|T.y| \leq AABB1.E.y + AABB2.E.y \ \&\&$$

$$|T.z| \leq AABB1.E.z + AABB2.E.z$$

### B) OBB

OBB (Oriented Bounding Box) sunt folosite pentru a încadra obiecte care sunt rotite și pentru care folosirea AABB conduce la determinarea unui volum încadrator prea mare, detecția coliziunii nefiind foarte precisă.

Un OBB este specificat prin centrul acestuia C, o mulțime de axe ortonormale  $A0, A1, A2$  și extensiile  $a0, a1, a2$  ( $a0 > 0, a1 > 0, a2 > 0$ ).

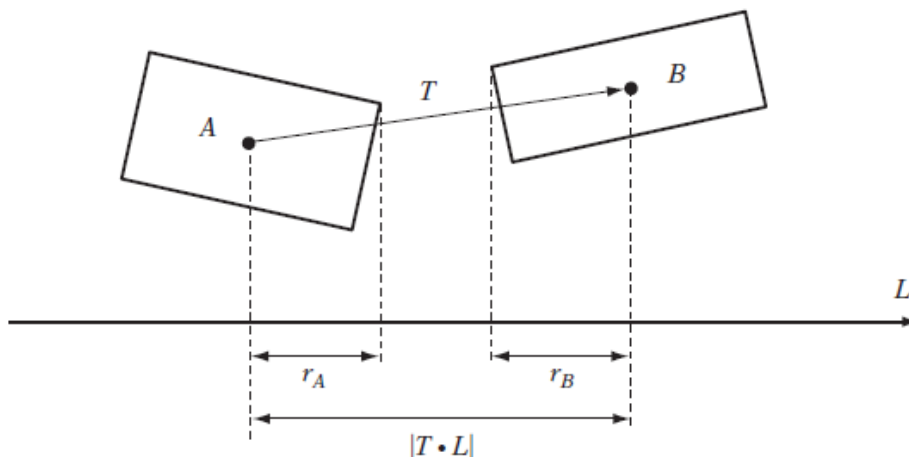
Un avantaj al OBB constă în faptul că acestea pot fi rotite odată cu obiectele, lucru care necesită însă reținerea orientării volumului.

Algoritmul folosit este foarte util deoarece încearcă minimizarea timpului necesar determinării intersecției a două obiecte.

Algoritmul folosește faptul că două poliedre convexe care nu se intersectează pot fi separate printr-un plan care, fie este paralel cu o față a unuia dintre poliedre, fie conține câte o muchie a fiecăruia dintre cele două poliedre.

Astfel, este necesar și suficient pentru a determina dacă două poliedre se intersectează sau nu, să se examineze intersecția proiecțiilor poliedrelor pe linii perpendiculare pe planurile de separație.





**Figura 4-18: Intersecție OBB**

Două OBB nu se intersectează atunci când suma extensiilor proiectate este mai mică decât distanța dintre proiecțiile centrelor:

$$|T \cdot L| > r_A + r_B$$

Pentru a determina corect starea de intersecție, este necesar să se facă verificări pentru un număr total de 15 axe:

- 3 pentru axele de referință ale primului OBB
- 3 pentru axele de referință al celui de-al doilea OBB
- 9 pentru axele perpendiculare pe planul format de oricare două din cele 6 axe de referință ale celor 2 OBB

Dacă exista cel puțin o axă pentru care nu se realizează intersecția, atunci obiectele nu se intersectează.

Numărul de operații poate fi redus prin exprimarea lui B în funcție de coordonatele lui A. Pentru aceasta se utilizează vectorul  $t$  de translație de la A la B și matricea de rotație  $R$  care aduce B în orientarea coordonatelor lui A.

### *C) Bounding-Sphere Collision Detection*

O modalitate destul de simplă de încadrare a unui obiect într-un volum o constituie determinarea unei sfere care să conțină obiectul respectiv.

Determinarea coliziunii a două sfere se poate face prin calculul distanței dintre centrele celor două sfere și compararea rezultatului cu suma razelor celor două sfere.

Astfel, dacă suma celor două raze este mai mică decât distanța calculată între centre, sferile nu se vor intersecta.

Dacă avem sferile reprezentate prin centrele  $C1(x_1, y_1, z_1)$  și  $C2(x_2, y_2, z_2)$  și razele  $r_1$  și  $r_2$ :

$$d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2 + (z_1 - z_2)^2}$$

,iar pentru intersecție este de ajuns să se compare  $d$  cu  $(r_1 + r_2)$

#### **4.3.2.2 Forțe implicate în simulările fizice**

Numărul de forțe simulate în aplicațiile grafice depind de gradul de realism pe care simularea dorește să îl obțină. Cu cât aplicația ține cont de mai multe forțe din mecanică, cu atât aceasta se apropie mai mult de realitate. Cele mai importante forțe de care trebuie să se țină cont sunt :

##### *A) Forța de frecare*

Forța de frecare reprezintă rezistența venită din partea mediului la mișcarea obiectului. Aceasta poate fi indusă de aer, un fluid sau de frecarea cu alte obiecte din scenă (la coliziuni). Forța de frecare decelerează mișcarea corpului cu un anumit factor, proporțional cu un indice de frecare corespunzător materialului obiectului.

Forța de frecare respectă formula:

$$F = \mu \cdot N$$

unde  $N$  este forța de apăsare normală, iar  $\mu$  este coeficientul de frecare

Exista 2 forme ale forței de frecare:

- frecare statică
- frecare cinetică

Frecarea statică apare atunci când cele 2 obiecte nu se află în mișcare relativă, deoarece forța de mișcare este mai mare decât forța de tracțiune.

Frecarea cinetică apare atunci când cele 2 obiecte se află în mișcare relativă. Coeficientul de frecare este în acest caz notat  $\mu_k$  și este în general mai mic decât coeficientul de frecare statică.

### *B) Forța de gravitație*

Prin folosirea acestei forțe se obține efectul de “cădere” al obiectelor lăsate liber în spațiu. Constă în accelerația obiectului (de obicei) în direcția negativă a axei verticale.

Forța de gravitație respectă formula:

$$G = g \cdot m$$

unde  $m$  este masa obiectului și  $g$  este constanta de accelerație gravitațională, aproximativ egală cu 9.8

Un corp având inițial viteza 0, lăsat liber, va accelera cu aproximativ  $9.8 \text{ m/s}^2$

În cazul general, viteza corpului la timpul  $t$ ,  $v(t)$  este:

$$v(t) = v(0) + g \cdot t$$

Distanța parcursă de corp până la timpul  $t$ ,  $d(t)$  este:

$$d(t) = v(0) \cdot t + g/2 \cdot t^2$$

Pentru a obține efecte realiste, forța de atracție gravitațională trebuie folosită împreună cu forța de frecare.

### **4.3.2.3 Algoritmi**

Au fost implementate următoarele calcule de fizică atât pe GPU cât și pe CPU (pentru comparație) :

- Detecția coliziunilor între obiecte
- Procesarea forțelor de ciocnire în urma coliziunilor
- Calcule pentru forța de gravitație

De asemenea, au fost folosite următoarele tipuri de obiecte asupra cărora au fost efectuate calculele de coliziune :

- Cuburi
- Sfere

- Cilindri

Aceste tipuri de obiecte reprezintă aproximările cele mai folosite în practică ca volume încadratoare a obiectelor 3D complexe pentru efectuarea calculelor de coliziune.

Algoritmul pentru lansarea calculelor de coliziune pe firele de execuție CUDA este următorul :

- Datele inițiale ale scenei sunt făcute disponibile în memoria globală a GPU-ului
- Setul de sarcini de coliziune este determinat de un vector unidimensional de dimensiune  $N*N$  (unde  $N$  este numărul de obiecte din scenă) deoarece se simulează coliziuni  $N$ -la- $N$  între obiecte
- Acest vector se împarte în sub-vectori de dimensiune fixă ce au fiecare dimensiunea 32
- Fiecare sub-vector va fi procesat de un bloc de fire de execuție unde fiecare fir de execuție va fi responsabil de calculele executate pentru un element al sub-vectorului (obiectul identificat de firul de execuție va fi testat pentru coliziuni cu toate celelalte obiecte ale scenei)
- Pentru implementarea *multi-GPGPU*, vectorul va fi împărțit la numărul de GPU-uri obținându-se astfel pentru fiecare GPU o porțiune de care acesta este responsabil (de exemplu dacă se folosesc două GPU-uri, GPU0 va procesa prima jumătate a vectorului și GPU1 va procesa cea de-a doua jumătate a vectorului)

Pentru implementarea *single-GPGPU* pseudocodul algoritmului este următorul :

```
Dimensiune bloc fire de execuție = (32,1)
Dimensiune grid fire de execuție =
( numarObiecteScena / dimensiuneXbloc, 1)

Copiază date în memoria globală GPU de pe mașina gazdă

Lansează în paralel pe GPU funcție kernel(dimensiune
bloc,dimensiune grid)

Sincronizează prin barieră terminarea firelor de execuție

Copiază date din memoria globală GPU pe mașina gazdă
```

Funcția kernel ce se execută pe fiecare fir de execuție de pe GPU este :

```
IndexObiect = IdBlocX * DimensiuneBlocX + IdFirDeExecuțieX

Pentru fiecare obiectJ din scena
    CalculeazăIntersectie( obiectDinScena[IndexObiect],
                           obiectDinScena[obiectJ] )
```

Pentru implementarea *multi-GPGPU* pseudocodul algoritmului se modifică astfel :

```
Dimensiune bloc fire de execuție = (32 , 1)
NrObiectePerGPU = numarObiecteScena / nrGPU
Dimensiune grid fire de execuție =
( NrObiectePerGPU / dimensiuneXbloc , 1)

Pentru fiecare GPU
{
    Copiază date în memoria globală GPU de pe mașina gazdă

    Lansează în paralel pe GPU funcție kernel(dimensiune
    bloc,dimensiune grid, idGPU)

    Sincronizează prin barieră terminarea firelor de execuție

    Copiază date din memoria globală GPU pe mașina gazdă
}
```

Funcția kernel ce se execută pe fiecare fir de execuție de pe GPU este :

```
Offset = NrObiectePerGPU * idGPU
IndexObiect = IdBlocX * DimensiuneBlocX + IdFirDeExecuțieX +
Offset

Pentru fiecare obiectJ din scena
    CalculeazăIntersectie( obiectDinScena[IndexObiect],
                           obiectDinScena[obiectJ] )
```

### 4.3.3 ANALIZĂ TIMP DE EXECUȚIE

Timpul de execuție pentru efectuarea calculelor de coliziune în implementarea pe CPU este următorul :

$$T(NrObiecte * NrObiecte ),$$

unde,

NrObiecte – numărul de obiecte din scenă

Timpul de execuție pentru efectuarea calculelor de coliziune în implementarea single GPGPU este următorul :

$$T( (NrObiecte * NrObiecte) / (nrSM * nrSP) )$$

+ *constantaOverhead*,

unde,

NrObiecte – numărul de obiecte din scenă

nrSM – numărul de multiprocesoare de pe GPU

nrSP – numărul de procesoare scalare de pe fiecare multiprocesor

*constantaOverhead* – constantă introdusă de operațiile de transfer de memorie și sincronizare

Timpul de execuție pentru efectuarea calculelor de coliziune în implementarea multi GPGPU este următorul :

$$T( (NrObiecte * NrObiecte) / (nrSM * nrSP * nrGPU) )$$

+ *nrGPU\*constantaOverhead*,

unde,

NrObiecte – numărul de obiecte din scenă

nrSM – numărul de multiprocesoare de pe GPU

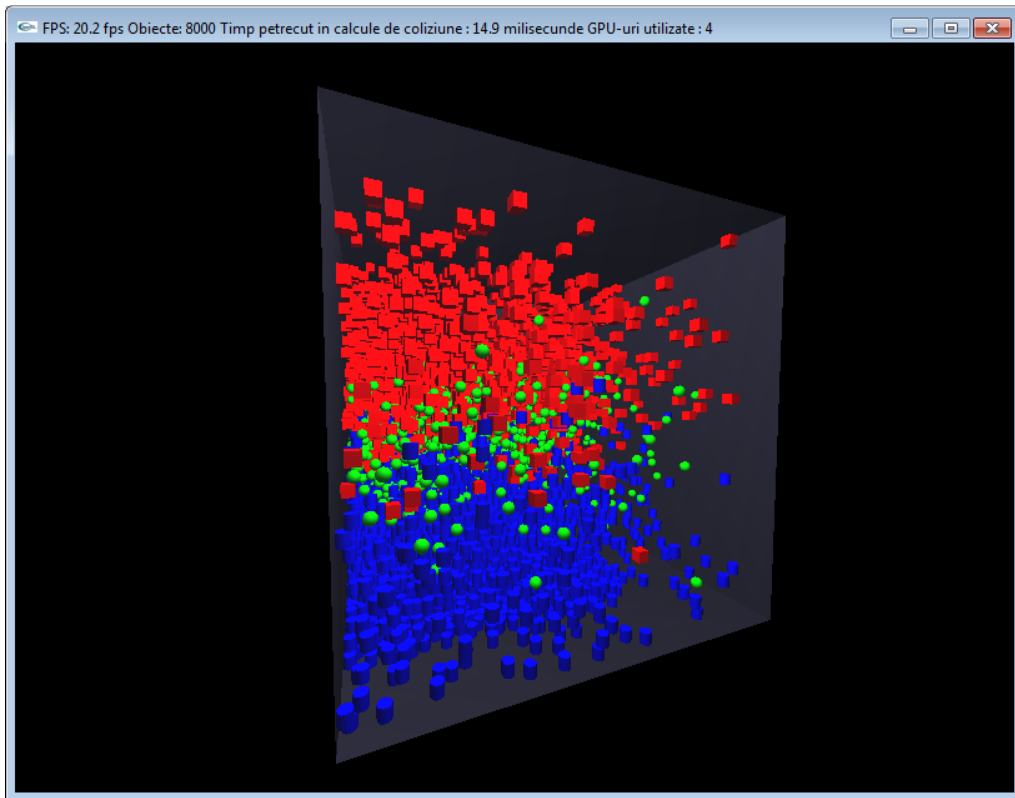
nrSP – numărul de procesoare scalare de pe fiecare multiprocesor

nrGPU – numărul de GPU-uri disponibile

*constantaOverhead* – constantă introdusă de operațiile de transfer de memorie și sincronizare

Testele au fost efectuate pe următoarele 3 implementări :

- CPU : Intel Core2Quad Q6600
- Single GPGPU : 1 core al unui GPU GTX590
- Multi GPGPU : 4 core-uri GPU (2 GPU GTX590)



*Figura 4-19: Simularea coliziunii a 8000 de obiecte folosind 4 core-uri GPU*



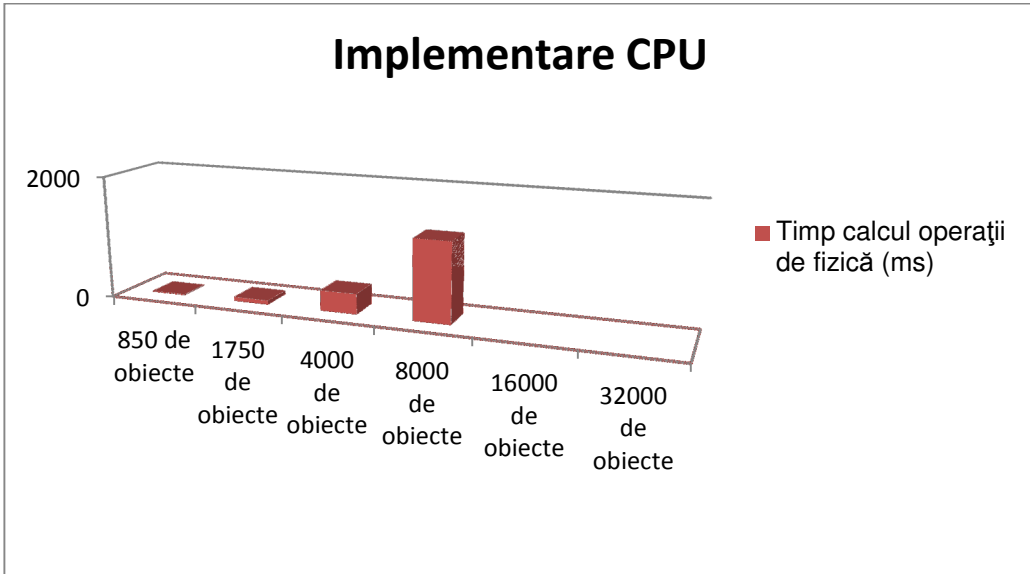
Rezultatele obținute sunt cele din tabelele și graficele de mai jos.

	850 de obiecte		1750 de obiecte		4000 de obiecte	
	FPS	Timp calcule fizică / iterație	FPS	Timp calcule fizică / iterație	FPS	Timp calcule fizică / iterație
<b>1 CPU Quad Core Q6600</b>	130	16ms	24.7	67ms	2.9	332ms
<b>Single GPGPU 1 core GPU GTX590</b>	202	2ms	98	3.8ms	44.6	7ms
<b>Multi GPGPU 4 cores GPU GTX590</b>	200	2.3ms	96	4ms	44	7.2ms

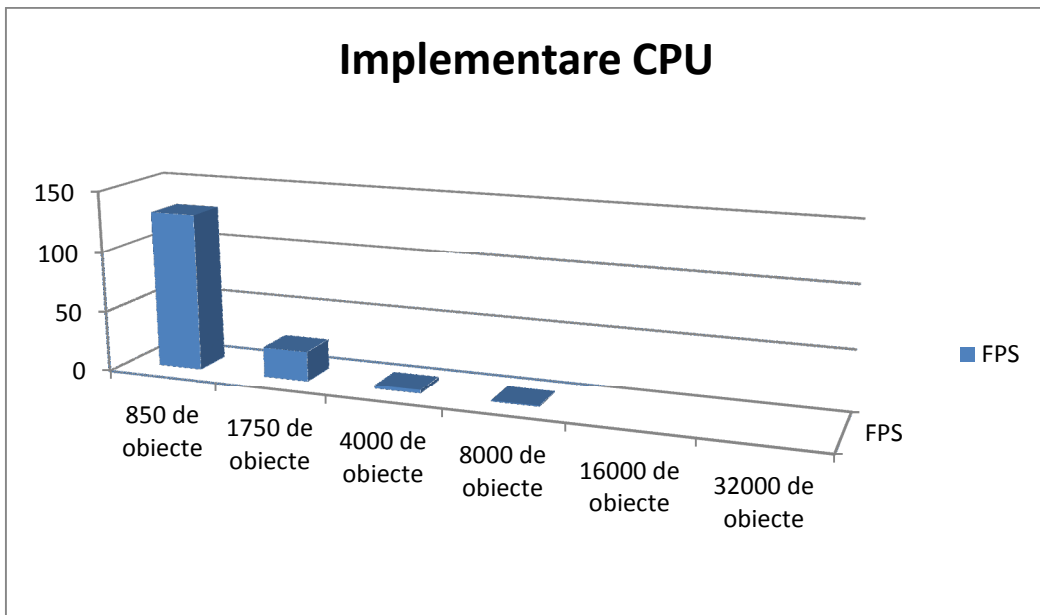
***Tabel 4-2: Rezultate implementare calcule fizică – Partea I***

	8000 de obiecte		16000 de obiecte		32000 de obiecte	
	FPS	Timp calcule fizică / iterație	FPS	Timp calcule fizică / iterație	FPS	Timp calcule fizică / iterație
<b>1 CPU Quad Core Q6600</b>	0.7	1314ms	N/A	N/A	N/A	N/A
<b>Single GPGPU 1 core GPU GTX590</b>	21.1	18ms	8	79ms	2.4	310ms
<b>Multi GPGPU 4 cores GPU GTX590</b>	22.4	15ms	11.6	30ms	5.4	80ms

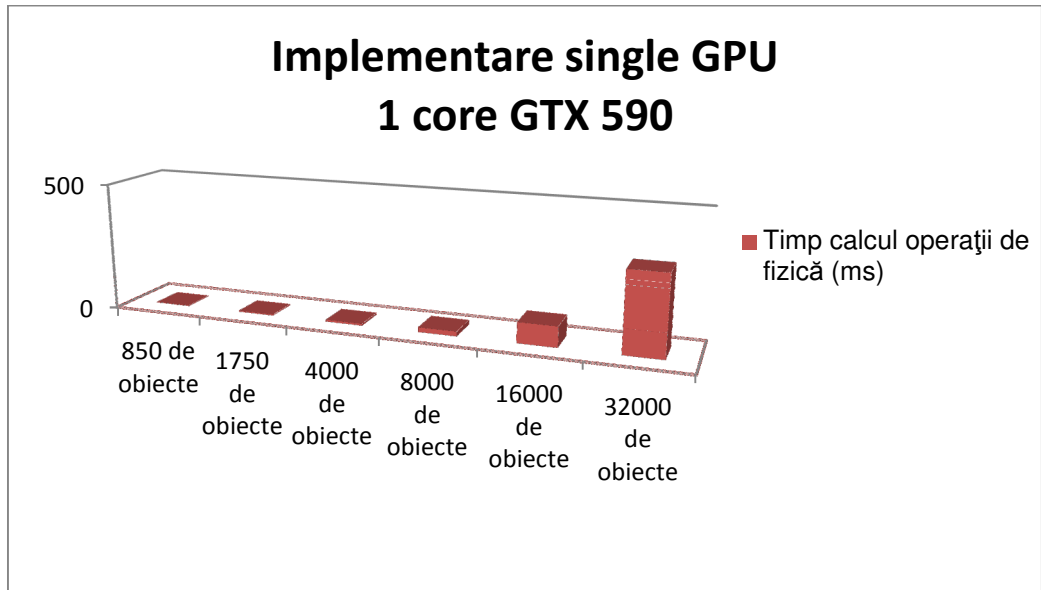
***Tabel 4-3: Rezultate implementare calcule fizică – Partea a II-a***



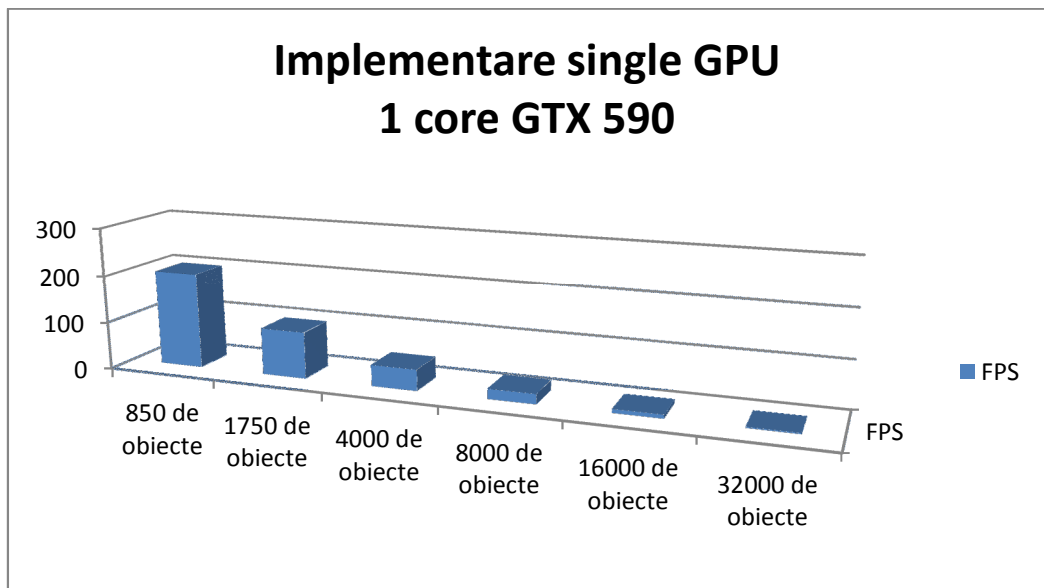
*Figura 4-20: Calcule de fizică pe CPU – Timpi de calcul*



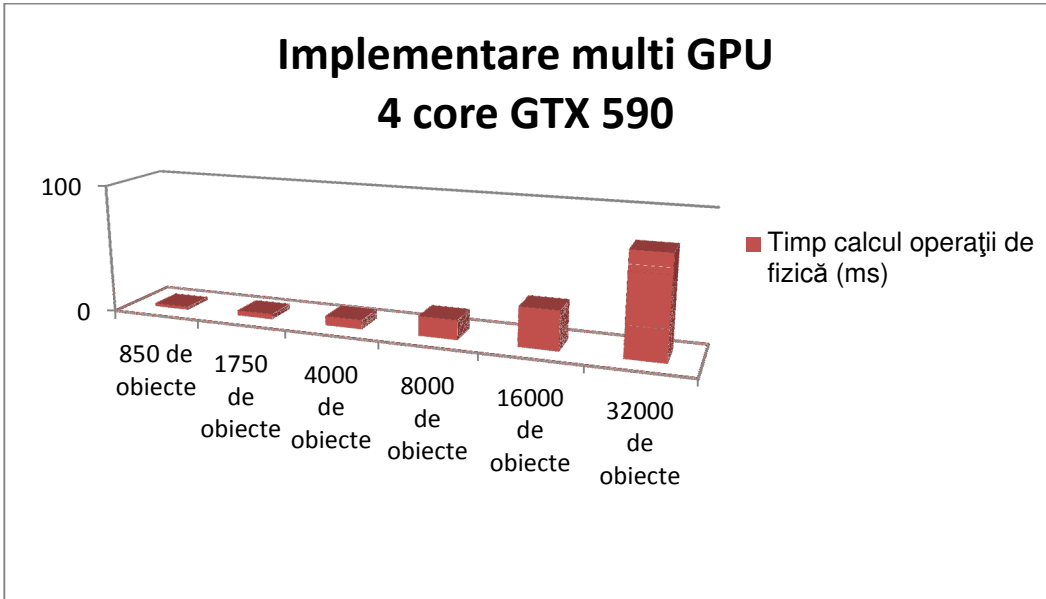
*Figura 4-21: Calcule de fizică pe CPU – FPS*



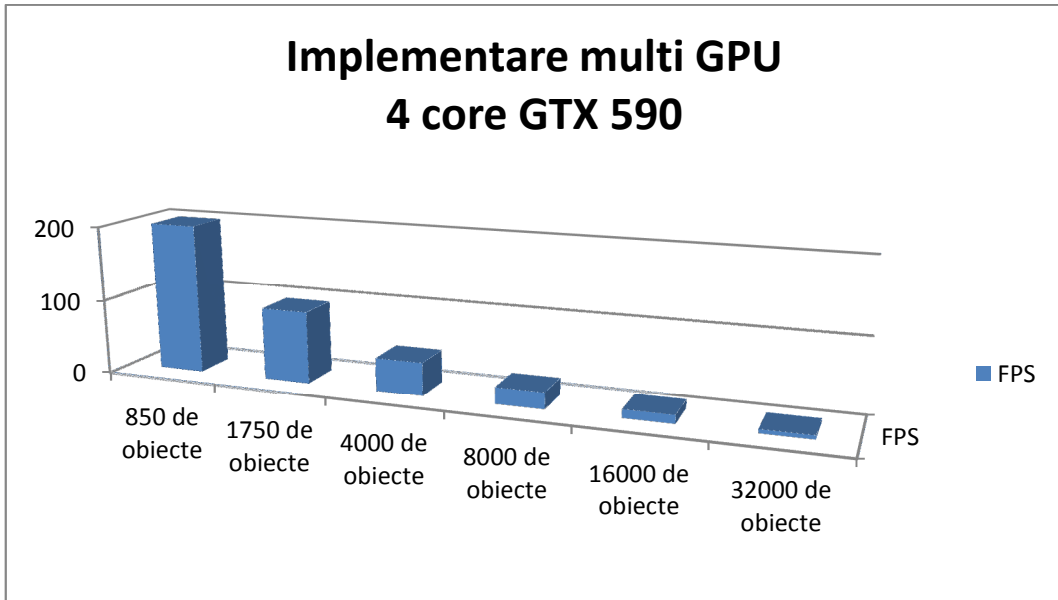
*Figura 4-22: Calcule de fizică pe 1 GPU – Timpi de calcul*



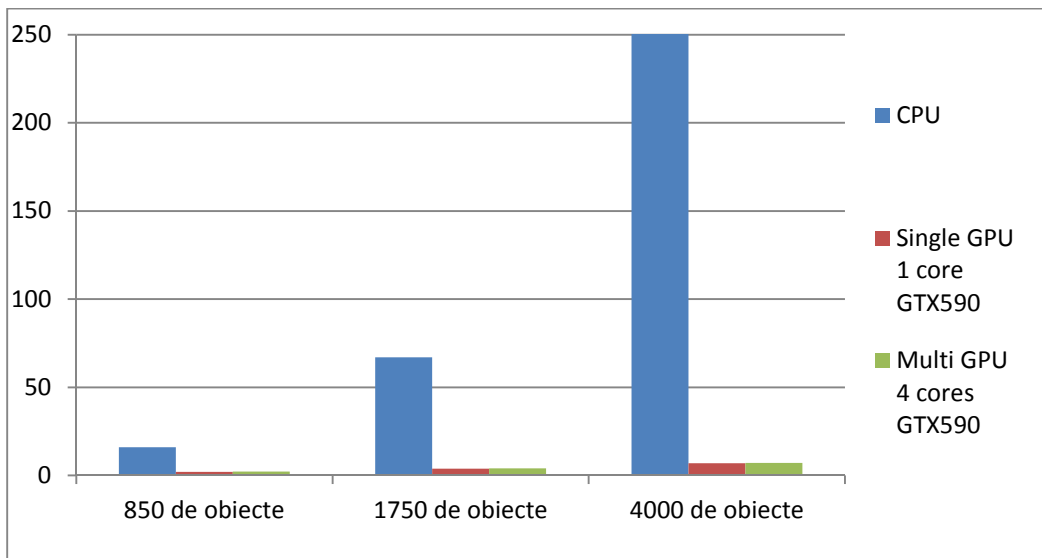
*Figura 4-23: Calcule de fizică pe 1 GPU – FPS*



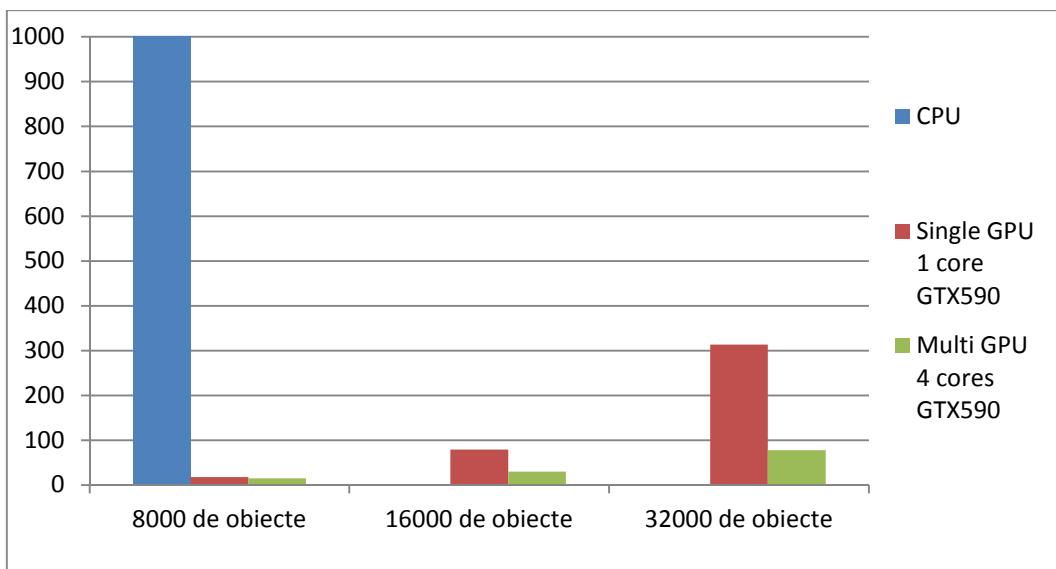
*Figura 4-24: Calcule de fizică pe 4 GPU – Timpi de calcul*



*Figura 4-25: Calcule de fizică pe 4 GPU – FPS*



*Figura 4-26: Comparație timp de calcul implementări operații de fizică executate în spațiile virtuale – Partea I*



*Figura 4-27: Comparație timp de calcul implementări operații de fizică executate în spațiile virtuale – Partea a II-a*

Din analiza rezultatelor de mai sus se pot trage următoarele concluzii :

- Pe implementarea CPU se poate observa faptul că timpul necesar pentru calculele de fizică crește rapid pe măsură ce crește și numărul de obiecte ce sunt implicate în simulare; pentru 4000 de obiecte în scenă, CPU deja nu mai face față, costul de timp fiind peste 300 de ms / iterație; de asemenea pentru 16000, respectiv 32000 de obiecte în scenă nu au putut fi obținute rezultate deoarece numărul de obiecte este prea mare pentru a putea fi procesat într-un timp rezonabil de CPU; având în vedere rezultatele de mai sus se poate deduce faptul că procesarea pe CPU a calculelor de fizică nu este scalabilă cu numărul de obiecte din scenă.
- Implementările GPGPU (single și multi) scalează bine o dată cu creșterea numărului de obiecte din scenă, comparativ cu implementarea pe CPU realizându-se sporuri de performanță semnificative; astfel pentru 1750 de obiecte în scenă se obține un *speedup* (spor de performanță) de aproximativ 17x, pentru 4000 de obiecte în scenă se obține un speedup de aproximativ 47x iar pentru 8000 de obiecte în scenă se obține un speedup de aproximativ 90x; de asemenea, se poate observa faptul că inclusiv atunci când există un număr foarte mare de obiecte în scenă (16000 respectiv 32000) timpii obținuți sunt acceptabili.
- Între implementarea single GPU (care folosește un singur core GPU GTX590) și implementarea multi GPU (care folosește 4 core GPU GTX590) diferențele de spor de performanță sunt mici pentru un număr relativ mic de obiecte din scenă, mai mult, implementarea single GPU fiind chiar un pic mai performantă decât cea multi GPU; acest lucru se întâmplă din cauza overhead-ului din implementarea multi GPU ce apare în principal datorită operațiilor de sincronizare și de transfer de memorie (gază-dispozitiv GPU și dispozitiv GPU-gază) ce trebuie efectuate pentru toate dispozitivele GPU din sistem; adevăratul spor de performanță se vede atunci când numărul de obiecte din scenă crește semnificativ, astfel pentru 16000 de obiecte în scenă avem un speedup de 2.66x pentru implementarea multi GPU față de cea single GPU, iar pentru 32000 de obiecte din scenă avem un speedup de 3.87x pentru implementarea multi GPU față de cea single GPU; se poate observa că pentru un număr mare de obiecte, implementarea multi GPU ajunge aproape de *speedup-ul* maxim teoretic (4x) față de cea *single-GPU*

- FPS-ul obținut este puternic dependent de capacitățile de afișare grafică ale GPU-ului; astfel, pentru un număr mare de obiecte (mai mare de 4000) nu mai este direct dependent de timpii pentru calculele de fizică ci doar de numărul mare de obiecte ce trebuie afișate.

## **4.4 CONCLUZII**

GPU modern a devenit un hardware foarte avansat în domeniul arhitecturilor de calcul paralele many-core și se diferențiază față de arhitectura clasică CPU prin utilizarea la scară largă a multithreading-ului hardware și a procesării SIMD – Single Instruction Multiple Data. În acest fel, se acordă prioritate operațiilor ce se pot executa în paralel asupra unor cantități mari de date față de execuția operațiilor *single task* de latență scăzută.

GPGPU (General Purpose on Graphical Processing Units) reprezintă o metodă prin care unitatea de procesare grafică poate fi utilizată pentru a executa calcule/programe de uz general. Acest lucru a devenit posibil odată cu apariția etapelor programabile ale GPU precum și datorită existenței bibliotecilor și instrumentelor de dezvoltare puse la dispoziție de producătorii de unități de procesare grafică.

Evoluția puterii de calcul a GPU moderne a înregistrat în ultima perioadă creșteri spectaculoase și sunt demne de menționat următoarele aspecte :

- viteza de dezvoltare a arhitecturilor GPU este una ridicată, de obicei la fiecare 12-18 luni fiind disponibile îmbunătățiri spre deosebire de ciclul de dezvoltare al CPU unde modificările la nivel arhitectural au loc la aproximativ 2 ani;
- soluțiile curente high-end de *procesoare* GPU deja au o putere de calcul de câteva zeci de ori mai mare decât cea a CPU atunci când sunt folosite pentru calcule asupra datelor ce se pretează la procesare SIMD.

Pentru a crea un grad de imersiune cât mai ridicat pentru utilizatorii spațiului virtual, la nivelul aplicației client se poate folosi ca metodă de redare algoritmul RayTracing.

Până de curând, utilizarea RayTracing ca metodă de redare în timp real nu era viabilă deoarece costul computațional al metodei este unul ridicat iar hardware-ul existent nu era suficient de rapid pentru a obține rezultate în timp real.

O dată cu lansarea la finalul anului 2006 a arhitecturilor unificate pentru GPU, a apărut posibilitatea pentru o gamă largă de operații de uz general să poată fi implementate folosind tehnici GPGPU.

În acest capitol au fost propuse soluții pentru a executa algoritmul RayTracing folosind procesare *single GPGPU* și *multi-GPGPU* și au fost prezentate rezultatele obținute. Este important de menționat că la timpul realizării cercetării (2008 și 2009) încercările de exploatare a paralelismului de tip GPGPU pentru implementarea algoritmului RayTracing se aflau într-o stare incipientă, astfel încât cercetarea realizată în acel moment a fost una de actualitate.

Soluțiile și rezultatele obținute au fost publicate într-o serie de articole ([*AMM09a*], [*AMM09c*]).

La nivelul serverelor de spații virtuale 3D MMO se execută un număr foarte mare de operații pentru a simula o fizică realistă a spațiului virtual (detectia coliziunilor, forțe de frecare, forța de gravitație, etc).

Acestea reprezintă gama de operații care sunt cele mai frecvent efectuate de către serverele spațiului virtual și astfel sunt principalele consumatoare de timp de calcul.

În acest capitol a fost propusă o soluție de paralelizare folosind procesare *single GPGPU* și *multi GPGPU* a calculelor de detecție a coliziunilor ce sunt executate în mod obișnuit în cadrul spațiilor virtuale.

Rezultatele obținute au scos în evidență următoarele aspecte majore :

- implementarea folosind CPU nu mai face față în timp real pentru un număr relativ mediu de obiecte existente în scenă; se poate deduce astfel că procesarea pe CPU a calculelor de fizică nu este scalabilă cu numărul de obiecte din scenă;



- implementarea ce folosește *single GPGPU* obține sporuri de performanță de ordinul zecilor față de CPU;
- implementarea ce folosește *multi GPGPU* obține sporuri de performanță de ordinul sutelor față de CPU;
- diferențele de spor de performanță între implementarea *single GPU* și implementarea *multi GPU* sunt mici pentru un număr relativ mic de obiecte din scenă, datorită overhead-ului introdus de operațiile de sincronizare și de transfer de memorie (gază-dispozitiv și dispozitiv-gază) ce trebuie efectuate pentru a gestiona mai multe GPU în sistem; adevăratul spor de performanță se vede când numărul de obiecte din scenă crește semnificativ, sporul de performanță crescând simțitor.

O parte din ideile utilizate pentru paralelizarea operațiilor de detecție a coliziunilor prin tehnici GPGPU au fost prezentate într-o serie de articole ([*AMM10a*] , [*AMM11*]).

## **5 ARHITECTURĂ ORIGINALĂ PENTRU SERVERE 3D MMO UTILIZÂND PARALELISMUL GPGPU**

Una dintre principalele probleme cu care se confruntă dezvoltatorii de spații virtuale 3D MMO este cea a încărcării foarte mare la care trebuie să facă față lumea virtuală. Aceasta încărcare este dată atât de numărul de sarcini, care este direct proporțional cu numărul de utilizatori on-line, cât și de complexitatea acestor sarcini care depinde de particularitățile lumii virtuale simulate.

### **5.1 NECESITATEA MODIFICĂRII ARHITECTURII CLASICE**

Deși anumite cerințe necesare pentru a accesa o aplicație MMO în condiții optime, cum ar fi o lățime de bandă mare, latență scăzută, devin din ce în ce mai accesibile utilizatorului de rând, acest lucru nu este suficient pentru a rezolva problemele de scalabilitate cu care se confruntă aplicațiile MMO în încercarea lor de a suporta cât mai mulți utilizatori.

Aceste probleme de scalabilitate apar în mare măsură datorită necesității de a menține consistența lumii virtuale. Încălcarea cerințelor de consistență poate să ducă, de exemplu, la apariția de artefacte vizuale (de exemplu, intrarea unui utilizator într-un zid al unei clădiri) care nu au consecințe pe termen lung, dar pot cu ușurință să apară probleme mult mai serioase cum ar fi, de exemplu, pierderea sau duplicarea de obiecte în urma unei tranzacții între utilizatorii spațiului virtual. Astfel, pe lângă scăderea realismului, pot să apară probleme importante din punct de vedere al securității.

În mod tradițional, marea majoritate a aplicațiilor MMO sunt implementate folosind o arhitectură de tip client-server care oferă următoarele avantaje:

- Control centralizat
- Securitate crescută
- Implementare relativ simplă

Deși această arhitectură, care este foarte răspândită, este potrivită pentru multe tipuri de aplicații distribuite, introduce și o serie de dezavantaje cele mai importante fiind :

- 1) **Scalabilitatea** – Performanța sistemului central de servere care simulează lumea virtuală reprezintă o sugrumare ajungându-se astfel la o limită a numărului total de utilizatori care pot accesa simultan lumea virtuală.
- 2) **Redundanța** – Pentru a asigura faptul ca un server este suficient de puternic astfel încât sa facă față momentelor de utilizare maximă, este necesară asigurarea unui grad ridicat de redundanță hardware pentru echipamentul respectiv
- 3) **Fiabilitatea** – Acest tip de arhitectură nu are o fiabilitate ridicată deoarece serverele reprezintă posibile puncte de cădere
- 4) **Costul** – În general, dezvoltarea unei aplicații MMO durează între 2-3 ani și costul de dezvoltare pentru a lansa o aplicație MMO de calitate medie este unul ridicat. Un alt aspect de luat în considerație este faptul ca aplicațiile MMO trebuie sa simuleze lumi virtuale de întindere foarte mare și acest lucru face necesară existența unui număr ridicat de servere, uneori de ordinul sutelor, pentru a găzdui în mod complet mediul virtual. Se poate ajunge astfel ca dupa lansarea unei aplicații MMO, costurile de întreținere să ajunga chiar și până la 80% din totalul încasărilor.

Ținând astfel cont de aspectele menționate anterior, apare necesitatea explorării de noi soluții privind arhitecturile pentru servere de spații virtuale 3D MMO și de asemenea de a găsi modalități de optimizare a operațiilor efectuate în cadrul simulării unui lumi virtuale pentru a elimina sau măcar reduce o parte din problemele curente cu care se confruntă aplicațiile MMO.

## **5.2 PREZENTAREA MODIFICĂRILOR**

După cum s-a menționat, principala problemă a arhitecturilor curente pentru servere de spații virtuale 3D MMO este cea a scalabilității acestora atunci când trebuie să facă față la un număr mare de utilizatori simultan deoarece aceștia generează un volum ridicat de operații ce trebuie executate pe partea de server.

Modificările propuse în continuare pentru a fi aplicate în cadrul arhitecturilor serverelor spațiilor virtuale încearcă să adreseze această problemă din două direcții:

- introducerea în arhitectura serverelor a capacităților de procesare a sarcinilor de calcul folosind tehnici GPGPU;
- propunerea unui mecanism original de alocare a resurselor necesare execuției calculelor ținându-se cont de o serie de factori ce pot influența alocarea; deoarece această alocare va fi strâns legată de tehnica de zonare ce este prezentă în cadrul tuturor tipurilor de arhitecturi, ne vom referi în continuare la acest mecanism ca reprezentând modificări aduse „operațiilor de zonare” prezente în cadrul arhitecturilor serverelor de spații virtuale 3D MMO.

Cele două direcții abordate sunt strâns legate între ele, modificările propuse la mecanismul de alocare al resurselor fiind făcute plecându-se de la premiza că în cadrul arhitecturii vor exista capacități de procesare masiv paralelă folosindu-se tehnici GPGPU.

## **5.3 DETALIEREA OPERAȚIILOR DE ZONARE**

Zonarea este unul dintre mecanismele foarte importante ale unei arhitecturi de spații virtuale MMO. La nivelul mecanismului de zonare se execută foarte multe operații ce tin de acomodarea utilizatorilor în aria geografică a spațiului virtual deservită de o anumită zonă logică.

De asemenea, un număr important de operații se execută pentru a asigura persistența spațiului virtual, fiind necesară o sincronizare perfectă între zonele adiacente / cu zonele de tranziție, existând astfel un flux crescut de informații.

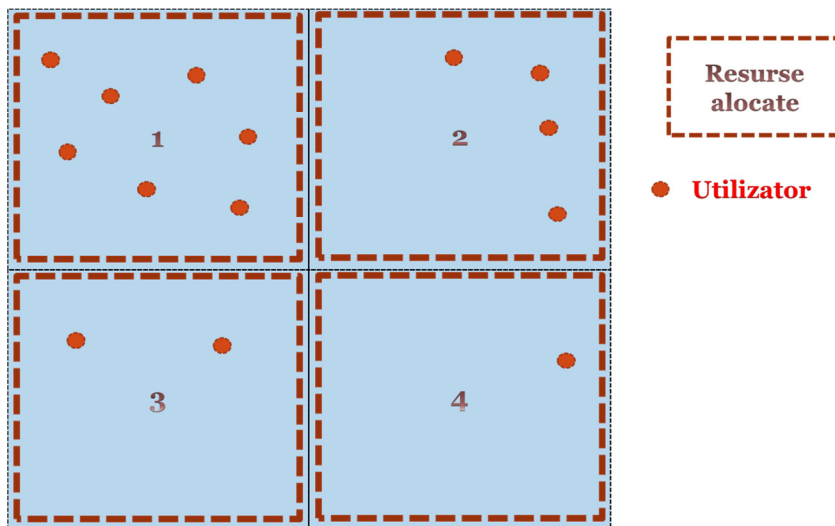
De exemplu, când un utilizator trece dintr-o zonă în alta care este adiacentă, este de dorit ca acest lucru să se facă într-un mod transparent pentru acesta și chiar dacă operațiile de transfer ce se efectuează între cele două zone sunt numeroase, acestea să nu fie vizibile într-un anumit fel utilizatorului (ecrane de încărcare, timp de răspuns, etc).

Așa cum s-a detaliat în capitolul 3, se folosesc în practică două abordări pentru a face zonarea unui spațiu virtual de o întindere geografică mare :

- cea a zonării simple, unde trecerea utilizatorului dintr-o zonă în alta este vizibilă acestuia (ecran de încărcare, obiectele dintr-o zonă nu sunt vizibile direct în zonele vecine, etc) realismul simulării având astfel de suferit;
- cea care realizează împărțirea pe zone logice într-un mod transparent pentru utilizator, prin utilizarea unor zone suplimentare speciale, denumite zone de tranziție, obținându-se astfel o simulare „seamless”, fără întreruperi cauzate de trecerea dintr-o zonă în alta, a spațiului virtual.

Indiferent de varianta de zonare folosită, există următorul dezavantaj prezent în implementările actuale : determinarea zonelor și a resurselor alocate pentru acestea nu se face dinamic, existând o configurație fixă sau care ține cont de foarte puțini parametri, aceste setări fiind stabilite încă de la începutul intrării în execuție a părții de server a aplicației MMO.

Astfel, indiferent de necesarul de resurse pentru o anumită zonă, există la nivel global o singură alocare bine stabilită pentru acea zonă.



*Figura 5-1: Alocare statică a resurselor*

Modificările prezentate în continuare pentru operațiile de zonare vizează optimizarea alocării resurselor, ținându-se cont de următorii 3 factori :

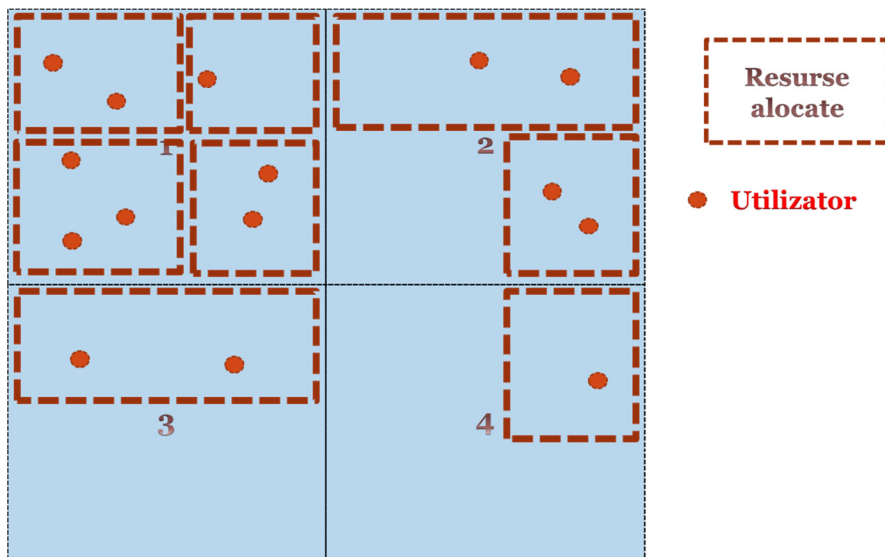
- poziționarea spațială a utilizatorilor care sunt on-line în cadrul spațiului;
- un cost de „greutate a operațiilor efectuate” asociat fiecărui utilizator care este on-line în cadrul spațiului virtual;
- tipurile de entități pentru care trebuie realizate calcule.

### 5.3.1 OPTIMIZARE CE ȚINE CONT DE POZIȚIONAREA SPAȚIALĂ ACTUALĂ A UTILIZATORILOR

O primă optimizare este dată de faptul că la alocarea resurselor se va ține cont de poziționarea spațială a utilizatorilor din cadrul spațiului virtual.

Astfel, dacă într-o zona geografică se află la un moment dat un anumit număr de utilizatori, alocarea resurselor pentru acea zonă poate fi ajustată dinamic (crescută sau scăzută) în funcție de numărul acestora.

În acest fel, nu se vor aloca resurse în mod inutil pentru zone unde nu există un număr mare de calcule de executat (din cauza lipsei sau a numărului mic de utilizatori ce pot genera operații de procesat), evitându-se astfel fenomenul de „înfometare” iar resursele disponibile vor putea fi utilizate în alte zone unde este nevoie de resurse suplimentare.



*Figura 5-2: Alocare dinamică a resurselor ce ține cont de poziționarea spațială a utilizatorilor*

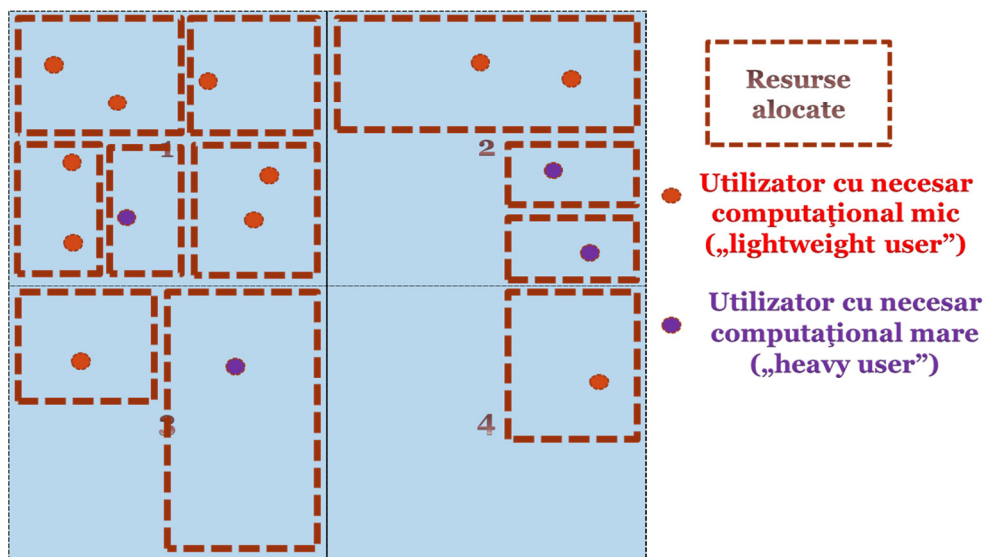
### 5.3.2 OPTIMIZARE CE ȚINE CONT DE COSTUL OPERAȚIILOR EFECTUATE DE FIECARE UTILIZATOR

O altă optimizare propusă ține cont de tipul de operații efectuate de un utilizator și de frecvența cu care acesta generează aceste operații.

De exemplu, dacă un utilizator nu se mișcă foarte des în spațiul virtual (este plecat din fața tastaturii, discută pe chat cu alți utilizatori, etc) el nu va genera operații ce sunt intensiv computaționale (cum ar fi cele de detecție a coliziunilor). Cazul opus celui prezentat anterior este acela în care utilizatorul execută foarte multe acțiuni ce presupun interacțiuni cu spațiul virtual (se află în tipul unui *quest*, lupte, etc) și atunci operațiile ce vor trebui executate de către server pentru gestiunea acestuia vor avea un cost computațional mare.

Prin execuția de operații de *profiling* în timp asupra acțiunilor unui utilizator din spațiul virtual, se poate astfel determina un cost de „greutate a operațiilor efectuate” asociat acestuia și în funcție de acest cost se poate schimba categoria din care face parte utilizatorul (utilizator care generează operații care nu sunt intensiv computaționale, utilizator care generează operații care sunt intensiv computaționale, etc).

În funcție de apartenența la diversele categorii existente se poate face o alocare a resurselor la un nivel mult mai fin, astfel încât să existe o mai bună utilizare a acestora.



*Figura 5-3: Alocare dinamică a resurselor ce ține cont de un cost asociat utilizatorilor*

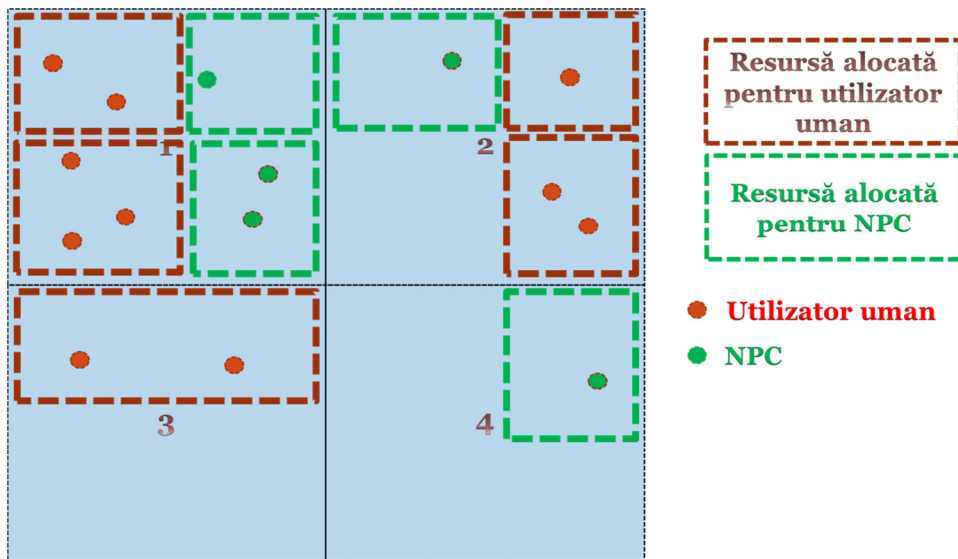
### 5.3.3 OPTIMIZARE CE ȚINE CONT DE TIPUL ENTITĂȚILOR PENTRU CARE SE EFECTUEAZĂ CALCULE

Ultima optimizare propusă vizează o alocare a resurselor care să țină cont și de tipul de entitate pentru care se vor efectua calcule. Mai exact, se va face o diferențiere între următoarele tipuri de entități ce populează spațiul virtual :

- utilizatori umani
- utilizatori controlați de server : NPC (non player characters)

Într-un spațiu virtual 3D MMO, cantitatea totală de NPC-uri raportată la numărul de utilizatori umani poate fi mai mare de zeci, chiar sute de ori.

De asemenea, NPC-urile fiind controlate de către modulul de inteligență artificială a serverului, au un alt comportament decât cel al unui utilizator normal astfel și operațiile generate de acestea sunt diferite față de cele ale unui utilizator uman.



*Figura 5-4: Alocare dinamică a resurselor ce țin cont de tipul entităților*

## 5.4 ABSTRACTIZAREA MATEMATICĂ A OPTIMIZĂRILOR PROPUSE

Pentru a putea folosi programatic în cadrul soluției arhitecturale optimizările propuse mai sus, se va folosi pentru determinarea unui cost asociat fiecărui task computațional declanșat de către un utilizator, următoarea funcție de cost :



$$\text{Cost}(\text{task}, \text{utilizator}) = \text{PozitionareSpatiala}(\text{task}, \text{utilizator}) * \text{pondereSpatiala} + \\ \text{ProfilingOperatii}(\text{task}, \text{utilizator}) * \text{pondereOperatii} + \\ \text{TipEntitate}(\text{task}, \text{utilizator}) * \text{pondereTipEntitate}$$

unde,

*PozitionareSpatiala(task,utilizator)*, *ProfilingOperatii(task,utilizator)* și *TipEntitate(task,utilizator)* sunt funcții asociate optimizărilor prezentate anterior și vor întoarce un cost determinat de acestea

iar,

*pondereSpatiala*, *pondereOperatii* și *pondereTipEntitate* sunt ponderi asociate costurilor determinate de cele 3 funcții care vor determina contribuția acestora la costul total al task-ului

Pseudocodul simplificat pentru funcția *PozitionareSpatiala(task,utilizator)* este următorul :

```
zonaUtilizator = determinaZonaGeografica(Utilizator)

numarUtilizatoriZona =
    determinaNumarUtilizatoriZonali(zonaUtilizator)

gradIncarcare =
    determinaGradIncarcareSpatial(zona, numarUtilizatoriZona)

întoarce gradIncarcare
```

Funcția *determinaGradIncarcare*, va întoarce valori asociate unor plaje discrete care sunt calculate în raport cu anumiți parametri ai spațiului virtual. De exemplu, ținând cont de numărul total al utilizatorilor din spațiul virtual și de valorile primite ca parametrii ai funcției (zona și numărul de utilizatori prezenți în aceasta) va întoarce valori asociate următoarelor grade de încărcare a zonei :

- încărcare mică
- încărcare medie
- încărcare mare

Pseudocodul simplificat pentru funcția *ProfilingOperatii (task,utilizator)* este următorul :

```
numarOperatiiUtilizator =  
    profiligNumarOperatii (Utilizator, interval de timp)  
  
tipOperatiiUtilizator =  
    profiligTipOperatii (Utilizator, interval de timp)  
  
gradIncarcare =  
    determinaGradIncarcareOperatii  
        ( Utilizator,  
          profiligNumarOperatii,  
          profiligTipOperatii )  
  
Întoarce gradIncarcare
```

Funcția *determinaGradIncarcareOperatii*, va întoarce valori asociate unor plaje discrete care sunt calculate în raport cu operațiile și tipul acestora executate de către utilizator într-un interval de timp. De exemplu, va întoarce valori asociate următoarelor categorii de utilizatori :

- utilizator cu necesar computațional redus
- utilizator cu necesar computațional mediu
- utilizator cu necesar computațional ridicat

Pseudocodul simplificat pentru funcția *TipEntitate (task,utilizator)* este următorul :

```
tipEntitate =  
    determinaTipEntitate (Utilizator)  
  
gradIncarcare =  
    determinaGradIncarcareEntitate(Utilizator,tipEntitate)  
  
Întoarce gradIncarcare
```

Funcția *determinaGradIncarcareEntitate*, va întoarce valori asociate tipului de entitate (uman sau NPC) a utilizatorului care a declanșat taskul. Astfel, pentru o

entitate de tip NPC, aceste valori pot depinde de o serie de parametri : numărul acestora comparativ cu cel al utilizatorilor umani, comportamentul acestora, etc.

Toți parametrii folosiți pentru determinarea diferitelor grade de încărcarea vor putea fi ajustați în cadrul implementării soluției arhitecturale pentru a realiza o calibrare cât mai fin posibilă a funcției globale de cost asociate unui task computațional.

## **5.5 PREZENTAREA SOLUTIEI ARHITECTURALE COMPLETE**

Pentru a putea introduce posibilitatea de execuție de operații GPGPU în cadrul arhitecturilor pentru servere de spații virtuale 3D MMO a fost necesară introducerea/modificarea și implementarea următoarelor componente :

- componenta de alocare a clienților și resurselor
- componenta pentru creare / planificare sarcini :
  - global la nivel de server
  - specializat la nivel de GPU
    - *single GPGPU*
    - *multi GPGPU*
- componenta pentru procesarea și propagarea rezultatului operațiilor CUDA

După cum se poate observa, modificările pentru a introduce capabilități GPGPU sunt în legătură strânsă cu modalitatea de alocare a clienților și resurselor prezentată anterior. Totodată, modificările aduse arhitecturii sunt gândite astfel încât să fie posibilă scalarea la nivel multi GPU.

Soluția prezentată în continuare are drept scop obținerea unui nivel ridicat de scalabilitate, atât vertical cât și orizontal, precum și reducerea limitărilor arhitecturilor tradiționale pentru servere de spații virtuale 3D MMO prin implementarea unei arhitecturi cu următoarele caracteristici importante :

- descarcă operațiile intensiv computaționale din sarcina CPU și le implementează ca operații GPGPU;

- este proiectată să fie masiv paralelă prin suportul planificării operațiilor la nivel de multi GPGPU;
- folosește un model care este condus de evenimente și sarcini, fiind centrat în principal pe locația, tipul și acțiunile executate de entitățile utilizator în cadrul spațiului virtual, deoarece ei sunt consumatorii puterii de calcul la nivel de server și nu datele spațiului virtual în sine.

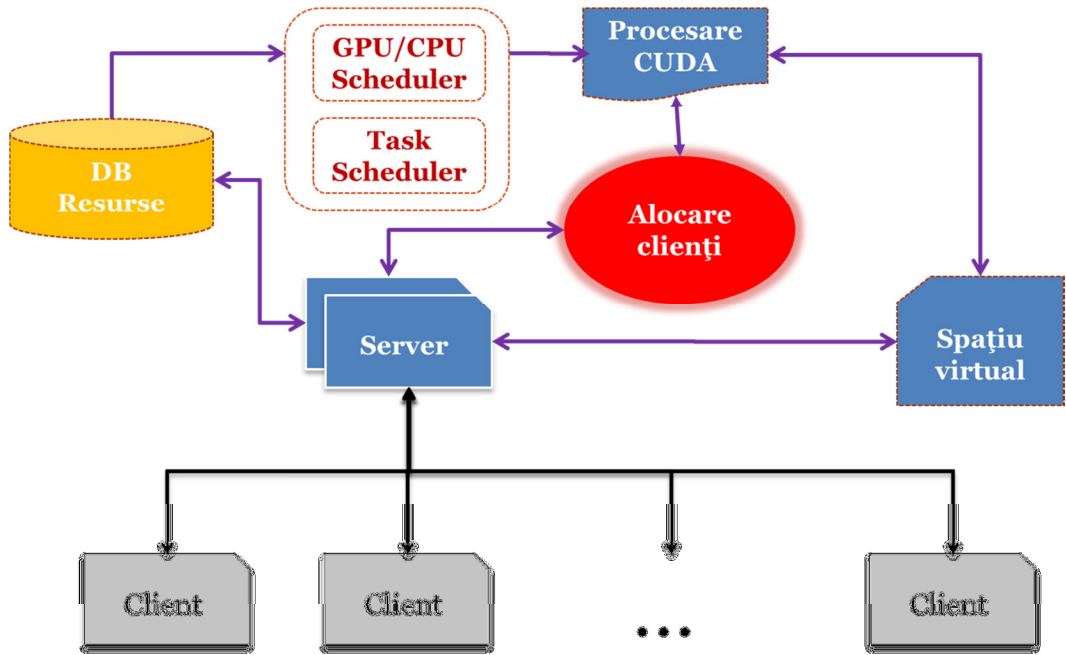
Pentru implementarea operațiilor GPGPU a fost folosit hardware grafic de la NVidia și toolkit-ul de dezvoltare CUDA (*[CUDAII]*). CUDA este o platformă hardware și software pentru executarea operațiilor de uz general, care în mod normal sunt executate de CPU, direct pe GPU fără a fi nevoie să se utilizeze un API grafic.

Următoarele operații vor fi implementate ca operații GPGPU :

- Operații pentru simularea fizicii spațiului virtual

Pentru a face posibilă execuția de operații GPGPU la nivelul arhitecturii serverelor de spații virtuale 3D MMO, următoarele module au fost implementate și utilizate la nivelul serverului :

- Modul de alocare clienți și resurse
- Modul de planificare sarcini (care are două componente distincte)
- Modul procesare CUDA



*Figura 5-5: Soluție arhitecturală cu posibilități de procesare GPGPU*

### **5.5.1 MODULUL ALOCARE CLIENȚI ȘI RESURSE**

Acest modul este responsabil cu alocarea clienților și a resurselor necesare executării calculelor la nivel de server.

Pentru realizarea alocării se folosesc următoarele criterii :

1. Poziționarea spațială în cadrul geografiei spațiului virtual a utilizatorilor care sunt on-line la momentul respectiv de timp. Astfel, dacă într-o anumită zonă logică sunt prezenți mai mulți utilizatori comparativ cu o altă zonă a spațiului virtual, celei dintâi îi vor fi alocate mai multe resurse de calcul din cele disponibile în sistem pentru a putea face astfel față la numărul mai mare de operații pe care aceasta le are de efectuat.
2. Tipul și numărul de operații executate de utilizatorii spațiului virtual. Astfel, alocarea procesării sarcinilor de calcul asociate clienților este realizată luându-se în considerație un cost de „greutate a operațiilor efectuate” asociat fiecărui utilizator care este on-line în cadrul spațiului virtual. În felul acesta, utilizatorii care generează un număr crescut de operații vor putea

beneficia de mai multe resurse pentru a se asigura faptul că sistemul va face față în timp real la acțiunile lor.

3. Tipul de utilizator pentru care trebuie realizate operații de calcul. Se face astfel o alocare diferită pentru utilizatorii de tip uman (care vor avea prioritate) și utilizatorii de tip NPC (non player characters)

Acest modul este responsabil de implementarea și utilizarea funcției globale a costului asociat unui task computațional, prezentate în paragraful 5.4. Astfel, se va aloca fiecărui task un cost ce va fi folosit mai departe în modulul de planificare a sarcinilor.

## **5.5.2 MODULUL PLANIFICARE SARCINI**

Acest modul este responsabil cu crearea și planificarea sarcinilor ce vor fi executate la nivelul serverelor spațiilor virtuale 3D MMO.

În arhitectura propusă, modulul de planificare a sarcinilor are două componente majore :

- componenta de creare a sarcinilor la nivel global (de sistem)
- componenta de planificare a sarcinilor la nivel de GPU / CPU

### **5.5.2.1 Creare sarcini la nivel de sistem**

Această componentă realizează crearea de sarcini de calcul asociate operațiilor executate în cadrul spațiului virtual. Astfel, la fiecare iterație (*tick*) a buclei principale a simulării, această componentă procesează acțiunile entităților spațiului virtual, creează sarcinile ce rezultă în urma procesării și le inserează într-o coadă de sarcini globală.

Tot această componentă este responsabilă și cu asigurarea menținerii ratei interne cu care funcționează spațiul virtual. De obicei, aceasta are o valoare ce asigura un FPS de 30 sau 60.

De asemenea, tot această componentă realizează trimiterea sarcinilor pentru a fi planificate la componenta de planificare.

### **5.5.2.2 Planificare sarcini la nivel de GPU / CPU**

Această componentă este responsabilă cu planificarea și trimiterea spre execuția propriu-zisă a sarcinilor la CPU sau la GPU.

Pentru implementarea GPGPU, folosind costul asociat fiecărui task, acest modul realizează planificarea sarcinilor ce vor fi executate pe multiprocesoarele GPU din sistem astfel :

- determină numărul de resurse de calcul (GPU) din sistem;
- în funcție de *tick-ul* intern și numărul de resurse de calcul, determină numărul maxim de sarcini ce pot fi executate la fiecare iterație a buclei principale a simulării spațiului virtual;
- în funcție de numărul de sarcini ce vor fi executate, determină dimensiunea blocurilor de fire de execuție și a grid-ului ce le conține;
- extrage din coada de sarcini globală acele sarcini ce vor fi planificate în *tick-ul* curent, în funcție de costul asociat acestora;
- creează vectorul de sarcini de executat folosind sarcinile extrase la pasul anterior și costul asociat acestora.

Pseudocodul simplificat pentru planificarea sarcinilor pentru calculul GPGPU este următorul :

```
nrGPU = determinaNumarResurseDeCalculDisponibile()

nrMaximSarcini = determinaSarciniDeExecutat
                (nrGPU, totalSarcini, tick intern)

dimBloc = determinaDimensiuneBlocFireDeExecutie()

nrSarciniPerGPU = nrMaximSarcini / nrGPU

dimensiune grid fire de execuție =
(nrSarciniPerGPU / dimensiuneXbloc , 1)

sarcini = extrageSarciniInFuncțieDeCost(listaGlobalaSarcini)

vectorSarcini = creeazaVectoSarcini(sarcini, nrMaximSarcini)
```

Un lucru important care a fost luat în considerație la proiectarea acestei componente a fost faptul ca aceasta să poată funcționa având la dispoziție mai multe resurse hardware pentru care să poată planifica sarcinile.

Astfel, s-a implementat fără mari dificultăți și soluția arhitecturală ce folosește procesare *multi-GPGPU*. În acest fel, se obține o scalabilitate transparentă în momentul în care în sistem sunt disponibile / introduse mai multe GPU.

### **5.5.3 MODULUL DE PROCESARE CUDA**

Acest modul primește sarcinile GPGPU de la modulul de planificare și este responsabil cu execuția acestora direct pe procesoarele plăcii grafice.

La nivelul acestui modul sunt realizate următoarele operații :

- implementarea funcțiilor kernel ce vor fi executate de firele de execuție ale procesoarelor scalare ale GPU-ului
- transferul de memorie de la dispozitivul gazdă la GPU
- după ce sarcinile au fost executate pe GPU, se realizează transferul în sens invers al memoriei (ce conține rezultatele execuțiilor) de la GPU la gazdă
- operațiile de sincronizare necesare pentru asigurarea faptului ca toate firele de execuție ale GPU-ului și-au terminat procesările



Pseudocodul simplificat pentru procesarea sarcinilor pentru calculul GPGPU este următorul :

```
Pentru fiecare GPU
```

```
{
```

```
    Copiază sarcini în memoria globală GPU de pe mașina gazdă  
    Copiază date spațiu virtual în memoria globală  
    GPU de pe mașina gazdă
```

```
    Lansează în paralel pe GPU funcție kernel(dimensiune  
    bloc,dimensiune grid, idGPU)
```

```
    Sincronizează prin barieră terminarea firelor de execuție
```

```
    Copiază rezultate sarcini din memoria globală GPU pe mașina  
    gazdă
```

```
    Copiază date spațiu virtual din memoria globală GPU pe  
    mașina gazdă
```

```
}
```

Tot acest modul este responsabil și cu propagarea rezultatelor execuțiilor către nivelul ce implementează logica spațiului virtual.

Folosind prototipul de arhitectură ce conține modulele prezentate mai sus s-au obținut rezultate încurajatoare în cadrul unor simulări. Acestea au permis interacțiunea în timp real a unui număr mare de utilizatori folosind soluții *single GPGPU* și *multi GPGPU*, bazate pe hardware grafic de la NVidia (GTX590).

Rezultatele obținute precum și o analiză a acestora vor fi prezentate în capitolul următor.

## **5.6 CONCLUZII**

Creșterea în popularitate a aplicațiilor MMO a dus la atragerea unui număr din ce în ce mai mare de utilizatori ce le accesează, acest fapt ducând la o încărcare foarte mare la care trebuie să facă față lumea virtuală în momentele de utilizare maximă. Aceasta încărcare este dată atât de numărul de sarcini, care este direct proporțional cu numărul de utilizatori on-line, cât și de complexitatea acestor sarcini care depinde de particularitățile spațiului virtual.

Problemele principale cu care se confruntă arhitecturile curente sunt :

- de scalabilitate
- de fiabilitate
- de redundanță
- de cost

Modificările propuse în prima parte a acestui capitol pentru a fi aplicate în cadrul arhitecturilor serverelor spațiilor virtuale încearcă să adreseze aceste probleme prin propunerea a două abordări originale :

- procesarea sarcinilor de calcul intensive la nivelul arhitecturilor folosind tehnici GPGPU;
- folosirea unui mecanism de alocare dinamică a resurselor necesare execuției calculelor și care este centrat pe utilizatorii spațiului virtual, deoarece aceștia generează necesarul computațional și nu datele lumii virtuale în sine.

În partea a doua a capitolului este propusă o soluție arhitecturală completă și sunt prezentate modulele necesare integrării în arhitectura serverelor de spații virtuale 3D MMO a abordărilor prezentate în prima parte a capitolului.

Soluțiile propuse și arhitectura ce le folosește au fost prezentate într-o serie de lucrări științifice ([AMB10], [MMA09b], [MMA09c]).

## **6 TESTAREA SCALABILITĂȚII SOLUȚIEI**

Pentru a realiza testarea de scalabilitate a fost realizat un prototip funcțional ce implementează funcționarea unui spațiu virtual 3D MMO bazat pe soluția arhitecturală propusă în capitolul anterior.

Astfel, au fost implementate module atât pentru partea de server cât și pentru cea de client.

### **6.1 SPECIFICAȚII DE FUNCȚIONARE SERVER**

La nivelul arhitecturii serverului au fost implementate modulele descrise în capitolul anterior :

- 1. Modulul de alocare a resurselor**
- 2. Modulul de creare și planificare a sarcinilor**
- 3. Modulul de procesare a sarcinilor**

Pentru a putea face o comparație între rezultate, toate cele 3 module au fost implementate astfel încât să funcționeze folosind procesare :

- CPU
- Single GPGPU
- Multi GPGPU

Serverul pune la dispoziție următoarele funcționalități către clienții care accesează spațiul virtual:

- Permite conectare clienți
- Autentificare clienți
- Trimitere evenimente către clienți :
  - o Actualizarea poziției acestora
  - o Actualizări legate de setările lumii virtuale (de exemplu, cu cât se poate deplasa la fiecare iterație în lumea virtuală)

- Actualizarea altor utilizatori conectați în spațiul virtual și care se află în aceeași zonă cu utilizatorul care a declanșat evenimentul de actualizare
- Actualizări legate de stări ale lumii virtuale
- Actualizări privind mesaje trimise de utilizatori între ei
- **Verificare corectitudine deplasare utilizator în spațiul virtual**
  - Utilizatorii se vor deplasa corect respectând următoarele restricții :
    - Nu vor putea să „intre” într-un alt utilizator
    - Nu vor putea să „intre” într-un alt obiect din spațiul virtual
    - Nu vor putea să „zboare”
    - Nu vor putea să ajungă sub teren
    - Nu se vor putea deplasa decât cu o viteză prestabilită; cu alte cuvinte dacă un utilizator va încerca să se deplaseze de la coordonatele curente la niște coordonate unde pentru a ajunge depășește viteza de deplasare impusă de spațiul virtual, mișcarea nu va fi considerată validă

Pentru simularea internă a logicii spațiului virtual la nivel de server au fost implementate următoarele funcționalități:

- Bucla principală a simulării spațiului virtual :
  - Nucleul unei aplicații MMO este reprezentat de o buclă de simulare care se execută sincronizat cu un *frame-rate*; pe parcursul fiecărei iterații ( *tick* ) al buclei de simulare, porțiunile ale stării sunt actualizate în concordanță cu logica simulării spațiului virtual; în implementarea de față a fost folosit un tick intern de 30 de milisecunde, echivalent cu un framerate de aproximativ 30 de cadre pe secundă;
- Descrierea geografiei lumii virtuale :
  - S-a folosit un teren simplu ce a fost generat dintr-o hartă de înălțimi. Aceasta geografie este „partajată” cu aplicația client, aceeași reprezentare fiind disponibilă fiecărui utilizator ce va accesa spațiul virtual.

Serverul pune la dispoziție către clienți următoarele funcții de API (Application programming interface) pentru comunicare :

- *Login()* : Conectare și autentificare la spațiul virtual

- *GetPosition()* : se obține poziția clientului în spațiul virtual; această poziție va fi folosită atunci când utilizatorul accesează de fiecare dată spațiul virtual pentru a fi „poziționat” unde trebuie, deoarece el poate părăsi spațiul virtual când se află în locuri diferite ale acestuia (serverul va stoca și pe disc pozițiile utilizatorilor astfel încât să poată să le furnizeze și dacă acesta este oprit / pornit din nou);
- *SetPosition()* : setează poziții pentru clienți ; aceasta funcție se va folosi în mod special de către server pentru a trimite unui utilizator al spațiului virtual, pozițiile celorlalți utilizatori (fie când aceștia se conectează, se deplasează, etc);
- *GetParameters()* : se întorc parametrii spațiului virtual (unul dintre ei foarte important, de care va trebui să se țină cont, este viteza de deplasare a utilizatorilor);
- *Move()* : utilizatorul încearcă deplasarea la o anumită poziție a spațiului virtual; i se va da voie sau nu (partea de decizie și de logică trebuie să fie implementată pe server pentru a preveni încercările de trișare ale clienților).

## **6.2 SPECIFICAȚII DE FUNCȚIONARE CLIENT**

Modulul client implementează următoarele funcționalități importante :

- interacțiunea cu utilizatorul;
- comunicarea cu serverul folosind interfața de programare (API) pusă la dispoziție de acesta;
- redarea 3D a spațiului virtual.

Modulul client oferă clienților ce accesează spațiul virtual următoarele:

- Conectare/Deconectare la spațiul virtual;
- Autentificare prin username/parolă;
- Interacțiuni cu utilizatorul
  - o Deplasarea acestuia în spațiul virtual
  - o Comunicare prin chat cu ceilalți utilizatori
- Clientul funcționează în următoarele 2 moduri
  - o Vizualizarea 3D a spațiului virtual :
    - Vizualizare a celorlalți utilizatori din spațiul virtual
    - Vizualizare a evenimentelor din spațiul virtual (deplasarea celorlalți utilizatori)

- Mod consolă :
  - În acest mod nu este necesară vizualizarea 3D a spațiului virtual; se vor primi în continuare toate evenimentele din spațiul virtual și se vor afișa informații legate de acestea în consolă; în continuare utilizatorul se va putea deplasa în spațiul virtual fie folosind taste direcționale, fie prin apelul direct din consolă al funcției de API *Move()*;
  - acest mod este necesar pentru a putea simula cu ușurință conectarea unui număr ridicat de utilizatori la spațiul virtual; astfel, se poate porni câte un proces/fir de execuție pentru fiecare client simulat, iar acesta poate primi un set de comenzi *Move()* pentru a avea un comportament în spațiul virtual.

### **6.3 COMUNICAREA ÎNTRE CLIENT ȘI SERVER**

Pentru a asigura posibilitatea procesării simultane a mai multor clienți, serverul utilizează o implementare bazată pe WinSockets. Astfel, aplicația server conține un fir de execuție pentru tratarea comunicației cu clienții și un alt fir de execuție pentru simularea lumii virtuale.

La pornirea unui client, acesta trimite serverului numele contului și parola asociate și, în caz ca sunt acceptate datele de autentificare, clientul primește de la server poziția și orientarea în lumea virtuală, precum și un identificator folosit în protocolul de comunicare.

Lumea virtuală este împărțită pe zone (în prototip, harta lumii virtuale are dimensiunea de 256x256 de unități și este împărțită în 16x16 zone). O regiune asociată unei zone reprezintă acea zonă și cele 8 zone înconjurătoare.

Comunicarea client-server este bazată pe evenimente. Astfel, clientul trimite mesaje serverului doar atunci când se realizează o schimbare de stare (execuția unei comenzi de mișcare din modul consolă, apăsarea sau ridicarea unei taste de mișcare din modul grafic, trimiterea unui mesaj de chat, ieșirea din aplicație).

Aceste evenimente au o arie de vizibilitate pentru client doar la nivelul regiunii în care se află. Astfel, la primirea unui mesaj ce indică declanșarea unui eveniment, serverul îl va propaga doar la clienții ce se află în regiunea clientului ce a generat evenimentul.

Atât clientul cât și serverul realizează o simulare a lumii virtuale cu datele care le au la dispoziție. De asemenea, testele de detecție de coliziune se realizează numai pe server.

Atunci când se detectează o coliziune, serverul trimite mesaje pentru a anunța acest eveniment tuturor clienților din regiunea în care s-a petrecut coliziunea. Acest mesaj conține pozițiile clienților (implicati în coliziune) dinainte de realizarea coliziunii.

În plus, serverul trimite și mesaje periodice tuturor clienților înregistrați în sistem, cu stările celorlalți clienți (poziție, orientare, stare de mișcare, destinație) din aceeași regiune. Acestea sunt folosite pentru sincronizarea clienților cu serverul. Prin acestea se elimină și necesitatea de a trimite mesaje cu starea clienților dintr-o nouă regiune atunci când clientul se mută în altă zonă.

Clientul conține și un timer asociat fiecărui client despre care el știe că se află în regiunea sa. Timer-ul asociat unui client este resetat atunci când se primește un mesaj ce face referire la acel client. Dacă un timer expiră, atunci se consideră ca acel client asociat timer-ului a părăsit regiunea (altfel s-ar fi primit informații despre el de la server prin mesajele periodice) și nu mai este nevoie să se păstreze informațiile despre el.

## **6.4 ALOCAREA RESURSELOR PENTRU CALCULUL COLIZIUNILOR**

Pentru a realiza calculul de coliziuni eficient se folosește o repartizare a clienților pe zone. În server se formează o reprezentare zonală a acestora, fiecare conținând toți clienții din zona asociată. Aceasta repartizare nu trebuie să se facă la fiecare iterație a buclei de simulare a lumii virtuale, deoarece clienții nu pot să străbată foarte repede o zonă.

Testele de coliziune sunt astfel limitate la nivelul unei zone și zonele adiacente (pentru a trata și cazul în care un avatar se află chiar pe marginea unei zone). În cadrul unui zone, testul de coliziune este de tip N-la-N, testându-se toate avatarele între ele.

Pentru modulul de detecție a coliziunilor pe GPU se folosește un grid în care fiecare bloc este asociat unei zone. În cadrul fiecărui bloc exista mai multe fire de execuție care realizează testele de coliziune de tip N-la-N, fiecare fir de execuție testând un anumit număr de avatare.

De asemenea, a fost implementat și mecanismul de alocare prezentat în capitoul anterior care ține cont de poziționarea spațială a utilizatorilor. Astfel, dacă un bloc de fire de execuție asociat unei zone nu are nimic de procesat (practic zona asociată nu conține utilizatori) acesta va fi realocat pentru a procesa teste de coliziune dintr-o zonă ce conține un număr mare de utilizatori, preluând astfel din efortul computațional necesar blocului care era asociat la zona țintă.

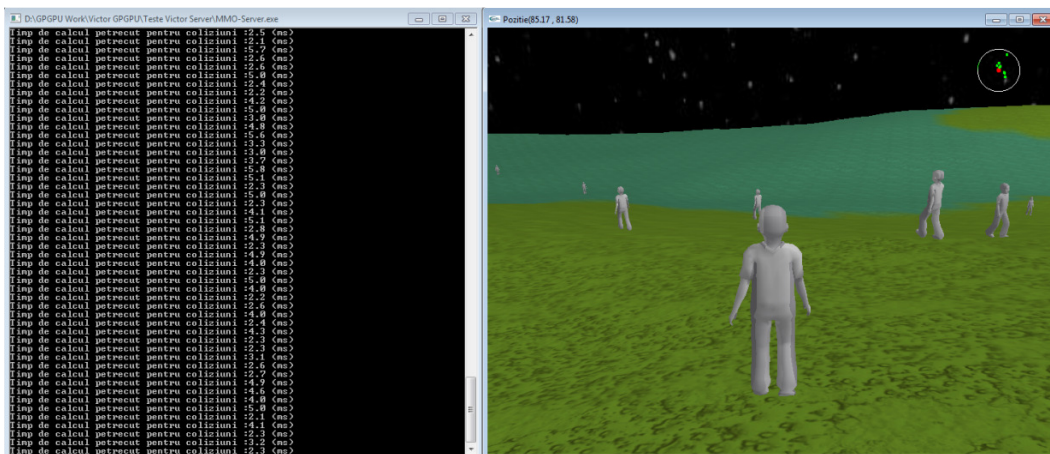


Figura 6-1: Prototip spațiu virtual ce implementează arhitectura propusă

## 6.5 REZULTATE

Pentru a simula încărcarea spațiului virtual cu un număr mare de clienți a fost folosită aplicația client în modul de funcționare consolă. Astfel, au fost pornite



pentru simularea clienților mai multe instanțe ale acesteia, fiecare primind un fișier ce conține comenzi de deplasare la poziții aleatoare din spațiul virtual cu scopul de a simula un comportament pentru fiecare client simulat.

Au fost realizate teste separate pentru următoarele numere de clienți ce au fost simulați că accesează spațiul virtual :

- 100 de clienți
- 1000 de clienți
- 4000 de clienți
- 8000 de clienți

Pentru efectuarea calculelor s-au testat toate cele 3 abordări :

- CPU : Intel Core2Quad Q6600
- Single GPGPU : 1 core al unui GPU GTX590
- Multi GPGPU : 4 core-uri GPU (2 GPU GTX590)

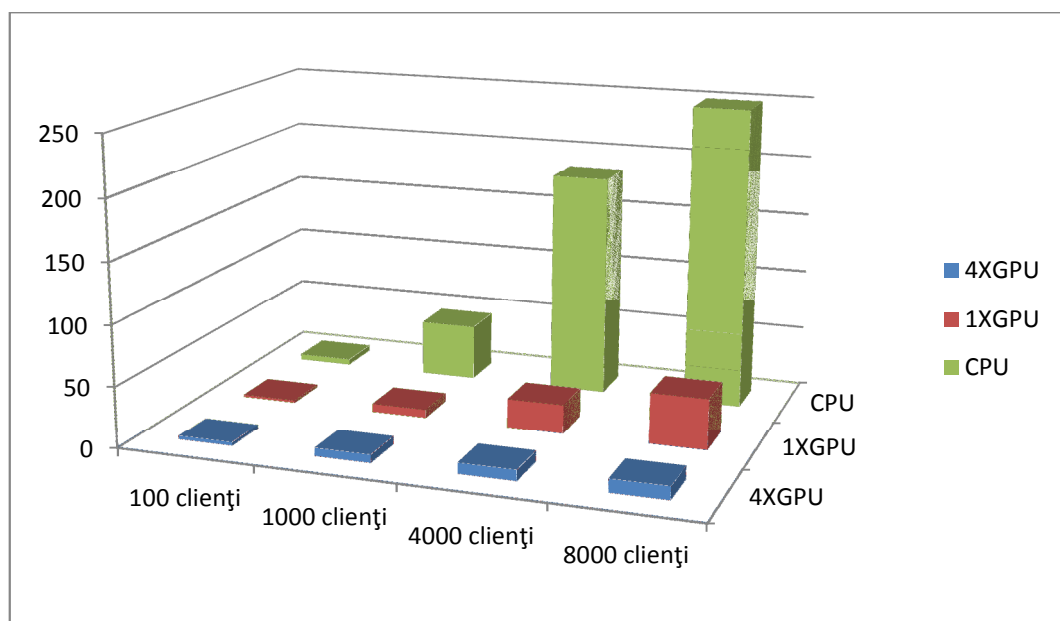
S-a avut în vedere timpul total necesar pentru calculul coliziunilor la fiecare iterație a buclei principale a spațiului virtual. *Tick-ul* intern folosit în implementarea serverului este de 30 de milisecunde, deci practic orice depășire a acestui prag conduce la imposibilitatea procesării în timp real pe partea de server a buclei principale de simulare, apărând întârzieri vizibile la nivelul aplicației client.

Rezultatele sunt prezentate în continuare.

## Testarea scalabilității soluției

	100 de clienți	1000 de clienți	4000 de clienți	8000 de clienți
	Timp calcule fizică / iterație	Timp calcule fizică / iterație	Timp calcule fizică / iterație	Timp calcule fizică / iterație
1 CPU Quad Core Q6600	4.7ms	46ms	185ms	365ms
Single GPGPU 1 core GPU GTX590	2ms	7ms	23ms	42ms
Multi GPGPU 4 cores GPU GTX590	3ms	6.8ms	9ms	11ms

**Tabel 6-1: Rezultate prototip funcțional al arhitecturii propuse**



**Figura 6-2: Comparație între timpii de execuție ai prototipului funcțional pentru cele trei implementari**

Din rezultatele obținute se pot astfel observa următoarele :

- pentru 1000 de clienți, implementarea pe CPU deja nu mai este în stare să facă față în timp real ratei interne la care funcționează spațiul virtual (30 de milisecunde);
- implementarea *single-GPU* are un speedup fata de cea CPU ce crește o dată cu numărul de utilizatori și de abia când numărul acestora atinge 8000 nu mai este în stare sa facă față la procesări în timp real;
- implementarea *multi-GPGPU*, datorită overhead-ului rezultat din operațiile de transfer de memorie și sincronizare, obține sporul de performanță așteptat față de cea *single-GPGPU* doar atunci când numărul de utilizatori este unul ridicat; astfel, chiar și pentru un număr de 8000 de utilizatori este capabilă să proceseze volumul de operații în timp real.

## **6.6 CONCLUZII**

În partea finală a tezei de doctorat s-a tratat testarea de scalabilitate a soluțiilor și arhitecturilor propuse, pentru a avea rezultate concrete de analizat.

În acest scop, a fost implementat un prototip funcțional al arhitecturii propuse, ce folosește :

- procesarea GPGPU a sarcinilor aferente operațiilor de detecție a coliziunilor;
- alocarea dinamică a resurselor GPU în funcție de poziționarea spațială a utilizatorilor în cadrul lumii virtuale.

În acest capitol au fost prezentate specificațiile de funcționare pentru aplicația server și client precum și modul în care s-a simulat un număr mare de clienți ce accesează spațiul virtual.

Rezultatele obținute au confirmat validitatea și viabilitatea soluției propuse.

## **7 CONCLUZII, CONTRIBUȚII ORIGINALE ȘI EVOLUȚII VIITOARE**

Popularitatea din ce în ce mai mare a spațiilor virtuale 3D MMO a condus inevitabil și spre creșterea semnificativă a numărului utilizatorilor care accesează astfel de spații. Acest lucru a condus la următoarele două consecințe majore :

- pentru a crea o experiență cât mai bogată și realistă pentru utilizatorii care accesează spațiile virtuale, aplicațiile MMO trebuie să poată oferi la nivelul aplicației client o redare 3D cât mai fidelă și realistă din punct de vedere grafic;
- a împins spre extrem încărcarea cu care se confruntă arhitectura acestor aplicații la nivel de server, uneori aceasta nemaifăcând față numărului foarte mare de utilizatori.

În prezent, arhitecturile de tip Client-Server pun la dispoziție funcționalitățile necesare unei aplicații MMO, însă o fac cu un cost ridicat. Avantajul unei soluții cu servere centralizate este dat de faptul că oferă un grad ridicat de control asupra sistemului, lucru care face posibilă asigurarea autentificării utilizatorilor, asigurarea persistenței și securității spațiului virtual. Astfel un cluster de servere performante poate oferi o putere computațională semnificativă iar centrele de date moderne pot face față unui trafic mare de date.

Din nefericire, natura puternic centralizată a arhitecturii Client-Server introduce o sugrumare din punct de vedere al performanței. În ciuda numărului lor care poate fi foarte mare, clusterurile de servere sunt de la un anumit punct limitate din punct de vedere computațional iar traficul de rețea este concentrat practic în echipamentele din centrele de date.

Costurile mari introduse de arhitectura Client-Server limitează practic scalabilitatea acesteia, fiind necesară împărțirea spațiului virtual în mai multe instanțe / *sharduri* (versiuni independente ale aceleiași lumi virtuale). De asemenea, deloc de neglijat este faptul că producătorii spațiilor virtuale 3D MMO trebuie să facă față unui cost financiar semnificativ pentru a menține funcțională infrastructura sistemului.

Din nefericire, arhitecturile Peer-To-Peer chiar dacă rezolvă câteva dintre neajunsurile arhitecturilor Client-Server, au la rândul lor câteva lipsuri semnificative care le fac de nefolosit la scară largă în aplicațiile MMO moderne.

Cel mai important neajuns este acela că nu oferă un mecanism viabil prin care să asigure persistența stării lumii virtuale. Atunci când un utilizator se deconectează de la spațiul virtual, resursele puse la dispoziție dispar, inclusiv datele pe care aplicația acestuia le gestiona.

Faptul că nu există un mecanism de autoritate centralizat face foarte dificilă actualizarea spațiului virtual de către producători și introduce posibile găuri de securitate.

De asemenea, deși participanții pot pune la dispoziție resurse suplimentare sistemului, acest lucru se realizează într-un mod eterogen, fiecare mașina având propriile limite computaționale. Fără un mecanism simplu de distribuire a încărcării, o concentrare mare de activități asupra unui peer poate să consume foarte rapid resursele acestuia.

Ținând cont de cele spuse mai sus, a fost realizată cercetarea posibilităților de îmbunătățire a tehnologiilor pentru spații virtuale MMO în următoarele direcții majore:

- la nivelul redării 3D oferite utilizatorilor:
  - posibilitatea creșterii realismului redării spațiului virtual 3D prin folosirea algoritmului RayTracing implementat utilizând tehnici de paralelizare GPGPU
- la nivelul serverelor de spații virtuale 3D MMO :
  - optimizarea operațiilor efectuate în cadrul spațiilor virtuale 3D MMO prin folosirea de tehnici GPGPU
  - modificări în arhitectura curentă a serverelor de spații virtuale 3D MMO

În continuare se va prezenta un rezumat al contribuțiilor originale așa cum reies din capitolele anterioare.

## **7.1 CONTRIBUȚII ORIGINALE**

### **➤ ANALIZA CONCEPTELOR, TEHNOLOGIILOR ȘI ARHITECTURILOR CURENTE FOLOSITE DE APLICAȚIILE MMO**

A fost efectuată o amplă analiză a conceptelor și tehnologiilor folosite pentru realizarea aplicațiilor MMO atât la nivel de client cât și de server.

De asemenea, au fost analizate arhitecturile curente folosite în implementarea spațiilor virtuale 3D MMO, punându-se în evidență atât avantajele utilizării acestora cât și lipsurile cu care se confruntă.

Această analiză, a dus la stabilirea unui punct de vedere propriu și original asupra tehnologiilor folosite în domeniu, stând la baza cercetărilor ulterioare și a soluțiilor propuse.

### ➤ **UTILIZAREA TEHNICILOR GPGPU PENTRU REDAREA 3D LA CLIENT**

Pentru a crea un grad de imersiune cât mai ridicat (prin creșterea realismului vizualizării lumii virtuale 3D) pentru utilizatorii spațiului MMO, la nivelul aplicației client se poate folosi ca metodă de redare algoritmul RayTracing.

Până de curând, utilizarea RayTracing ca metodă de redare în timp real nu era viabilă deoarece costul computațional al metodei este unul ridicat iar hardware-ul existent nu era suficient de rapid pentru a obține rezultate în timp real.

Odată cu lansarea la sfârșitul anului 2006 a arhitecturilor unificate pentru GPU, a apărut posibilitatea ca o gamă largă de operații de uz general să poată fi implementate folosind tehnici GPGPU.

Au fost astfel propuse soluții pentru a executa algoritmul RayTracing folosind procesare *single-GPGPU* și *multi-GPGPU*. Soluțiile propuse au fost implementate și au fost prezentate rezultatele obținute.

Este important de menționat că la timpul realizării cercetării (2008 și 2009) încercările de exploatare a paralelismului de tip GPGPU pentru implementarea algoritmului RayTracing se aflau într-o stare incipientă astfel încât cercetarea realizată în acel moment a fost una de actualitate.

### ➤ **UTILIZAREA TEHNICILOR GPGPU ÎN CADRUL SERVERELOR MMO**

La nivelul serverelor de spații virtuale 3D MMO se execută un număr foarte mare de operații pentru a simula o fizică realistă a spațiului virtual (detectia coliziunilor, forțe de frecare, forța de gravitație, etc).

Acestea reprezintă gama de operații care sunt cele mai frecvent efectuate de către serverele spațiului virtual și astfel sunt principalele consumatoare de timp de calcul.

A fost propusă o soluție de paralelizare folosind procesare *single-GPGPU* și *multi-GPGPU* a calculelor de detecție a coliziunilor ce sunt executate în mod obișnuit în cadrul spațiilor virtuale. Soluția a fost și implementată.

Rezultatele obținute în urma implementării soluției au confirmat validitatea ei și au scos în evidență următoarele aspecte majore :

- implementarea folosind CPU nu mai face față în timp real pentru un număr relativ mediu de obiecte existente în scenă; se poate deduce astfel că procesarea pe CPU a calculelor de fizică nu este scalabilă cu numărul de obiecte din scenă;
- implementarea ce folosește *single-GPGPU* obține sporuri de performanță de ordinul zecilor față de CPU;
- implementarea ce folosește *multi-GPGPU* obține sporuri de performanță de ordinul sutelor față de CPU.

### ➤ **PROPUNEREA UNUI MECANISM DE ALOCARE DINAMICĂ A RESURSELOR, CENTRAT PE UTILIZATORII SPAȚIULUI VIRTUAL**

A fost propusă folosirea unui mecanism de alocare dinamică a resurselor necesare execuției calculelor, care este centrat pe utilizatorii spațiului virtual deoarece aceștia generează necesarul computațional și nu datele lumii virtuale în sine.

Soluția propusă vizează optimizarea alocării resurselor ținându-se cont de următorii 3 factori :

- poziționarea spațială a utilizatorilor care sunt on-line în cadrul spațiului virtual;
- un cost de „greutate a operațiilor efectuate” asociat fiecărui utilizator care este on-line în cadrul spațiului virtual;
- tipurile de entități pentru care trebuie realizate calcule.

➤ **SOLUȚIE ARHITECTURALĂ COMPLETĂ CU POSIBILITATEA EXPLOATĂRII PARALELISMULUI DE TIP GPGPU**

A fost propusă o soluție arhitecturală completă ce conține modulele necesare integrării în arhitectura serverelor de spații virtuale 3D MMO a soluțiilor de optimizare prezentate anterior :

- capabilități de procesare a sarcinilor de calcul folosind tehnici GPGPU;
- folosirea mecanismului de alocare a resurselor centrat pe utilizatori.

➤ **IMPLEMENTAREA UNUI PROTOTIP FUNCȚIONAL PENTRU ARHITECTURA PROPUȘĂ ȘI TESTAREA DE SCALABILITATE A ACESTEIA**

A fost implementat un prototip funcțional al arhitecturii propuse, ce folosește:

- procesare GPGPU a sarcinilor aferente operațiilor de detecție a coliziunilor;
- alocarea dinamică a resurselor GPU în funcție de poziționarea spațială a acestora în cadrul lumii virtuale.

A fost implementată aplicația server, aplicația client precum și un mecanism de simulare a unui număr mare de clienți ce accesează spațiul virtual pentru a se putea realiza testele de scalabilitate.

Rezultatele obținute au confirmat validitatea și viabilitatea soluției arhitecturale propuse.

## **7.2 EVOLUȚII VIITORE**

În viitor se poate continua cercetarea în următoarele direcții, pentru a îmbunătăți / rafina soluțiile prezentate în această lucrare :



- implementarea la nivelul prototipului funcțional și a celorlalți factori propuși a fi folosiți de mecanismul de alocare :
  - tipul și frecvența operațiilor executate de utilizatorii spațiului virtual
  - tipurile de entități care generează sarcini de calcul
  
- ameliorarea overhead-ului introdus de operațiile de transfer de memorie și sincronizare efectuate atât pentru implementările *single-GPGPU* cât și *multi-GPGPU*; se pot cerceta soluții de optimizare a acestor operații consumatoare de timp folosind noile capacități de adresare / transfer a memoriei disponibile în arhitectura 4.0 CUDA ce a apărut în Mai 2011 :
  - posibilitatea de partaja accesul la GPU din contextul mai multor fire de execuție independente;
  - utilizarea tuturor GPU din sistem în mod concurent din cadrul unui singur fir de execuție;
  - **spațiu de memorie virtual unificat** : existența unui singur spațiu de adresare (CPU și GPU) pentru o utilizare mai facilă a resurselor de memorie;
  - **acces Peer-To-Peer** : posibilitatea de comunicare directă între Unitățile de Prelucrare Grafică (GPU) din sistem.

## **8 LISTA LUCRĂRILOR ȘI ACTIVITĂȚILOR ȘTIINȚIFICE ALE AUTORULUI**

### **8.1 LUCRĂRI LEGATE DE TEMATICA TEZEI**

#### **8.1.1 CĂRȚI**

[MMA09a] Alin Moldoveanu, Florica Moldoveanu, **Victor Asavei**, Costin Boiangiu  
Realitatea Virtuală  
Editura Matrix Rom, ISBN:978-973-755-488-8, 216 pg., 2009

#### **8.1.2 ARTICOLE**

[AMM09a] **Victor Asavei**, Florica Moldoveanu, Alin Moldoveanu  
Ray Tracing as GPGPU  
CSCS 17 - The 17th International Conference On Control Systems And Computer Science, 26-29  
May 2009, Bucharest

[AMM09b] **Victor Asavei**, Vlad-Valentin Ionita, Florica Moldoveanu, Alin Moldoveanu  
Real-Time Parallel Volume Rendering using Nvidia Cuda for Medical Imaging  
1483-1485, Annals of DAAAM for 2009 & Proceedings of the 20th International DAAAM Symposium,  
ISBN 978-3-901509-70-4, ISSN 1726-9679, pp 742, Editor B[ranko] Katalinic, Published by  
DAAAM International, Vienna, Austria 2009 (**Proceedings ISI**)

[AMM09c] **Victor Asavei**, Florica Moldoveanu, Alin Moldoveanu, Costin-Anton Boiangiu, Ruxandra  
Marasescu  
Ray Tracing as Multi GPGPU  
1203-1205, Annals of DAAAM for 2009 & Proceedings of the 20th International DAAAM Symposium,  
ISBN 978-3-901509-70-4, ISSN 1726-9679, pp 602, Editor B[ranko] Katalinic, Published by  
DAAAM International, Vienna, Austria 2009 (**Proceedings ISI**)

[AMM10a] **Victor Asavei**, Alin Moldoveanu, Florica Moldoveanu, Anca Morar, Alexandru Egner  
GPGPU for Cheaper 3D MMO Servers

International Conference on Telecommunications and Informatics : Session Information Science and Applications, Catania, Italia, 29-31 Mai 2010 (**Proceedings ISI**)

[AMB10] **Asavei, V.**; *Moldoveanu, A. D. B.; Moldoveanu, F.; Boiangiu, C. A.; Morar, A. - A. & Egner, A. I.*

**Innovative 3D MMO Servers Architectures Based on GPGPU**

Annals of DAAAM for 2010 & Proceedings of the 21st International DAAAM Symposium, ISBN 978-3-901509-73-5, ISSN 1726-9679, pp 0325, Editor B. Katalinic, Published by DAAAM International, Vienna, Austria 2010 (**Proceedings ISI**)

[AMM11] **Victor Asavei**, *Florica Moldoveanu, Alin Moldoveanu, Alexandru Egner, Anca Morar*

**Multi GPGPU Optimizations for 3D MMO Virtual Spaces**

CSCS 18 - The 18th International Conference On Control Systems And Computer Science, 24-27 May 2011, Bucharest

[MMA08a] *Alin Moldoveanu, Florica Moldoveanu, Alexandru Soceanu, **Victor Asavei***

**A 3D Virtual Museum**

Scientific Bulletin of UPB, Series C, vol. 70, No.3, ISSN 1454-234x, 2008

[MMA08b] *Alin Moldoveanu, Florica Moldoveanu, **Victor Asavei***

**Méga Musée Virtuel**

Conférence Internationale Francophone d'Automatique (CIFA) , ISBN: 978-2-915913-24-8 , Bucarest, 3-5 Septembre 2008

[MMA09b] *Alin Moldoveanu, Florica Moldoveanu, **Victor Asavei***

**Highly Scalable Server Architecture for Massive Multi-player 3D Virtual Spaces**

"Recent Advances in Applied Mathematics and Computational and Information Sciences", vol II, Proceedings of the 15th American Conference on Applied Mathematics and Proceedings of the International Conference on Computational and Information Science 2009, Houston, USA, April 30-May 2, 2009, ISSN: 1790-5117, ISBN: 978-960-474-071-0

[MMA09c] *Alin Moldoveanu, Florica Moldoveanu, **Victor Asavei***

**More scalability at lower costs – Server Architecture for Massive Multi-player 3D Virtual Spaces powered by GPGPU**

International Journal of Computers and Communication, Issue 2, Volume 1, 2009, ISSN: 2074-1294, published by University Press, London, UK, 2009

[MMA09d] *Alin Moldoveanu, Florica Moldoveanu, **Victor Asavei***

**Méga Musée Virtuel**

Automatique Avancee et Informatique Appliquee, Editura Academiei Romane, 2009

[MMA11] *Alin Dragos Bogdan Moldoveanu, Florica Moldoveanu, **Victor Asavei**, Alexandru Egner, Anca Morar*

From HTML to 3DMMO - a Roadmap Full of Challenges

CSCS 18 - The 18th International Conference On Control Systems And Computer Science, 24-27 May 2011, Bucharest

[LMM11] *Lucian Petrescu, Anca Morar, Florica Moldoveanu and **Victor Asavei***

Real Time Reconstruction of Volumes from Very Large Datasets Using CUDA

International Conference on System Theory, Control and Computing, 14-16 October, Sinaia, Romania (accepted as a full paper and for publication in the conference proceedings)

## **8.2 ALTE LUCRĂRI PUBLICATE ÎN PERIOADA DESFĂȘURĂRII DOCTORATULUI**

[MMM10a] *Anca Morar, Florica Moldoveanu, Alin Moldoveanu, **Victor Asavei**, Alexandru Egner*

Computer Assisted Analysis of Orthopedic Radiographic Images

International Conference on Signal Processing : Session Image Processing, Catania, Italia, 29-31 Mai 2010  
**(Proceedings ISI)**

[MMM10b] *Anca Morar, Florica Moldoveanu, Alin Moldoveanu, **Victor Asavei**, Alexandru Egner*

Medical Image Processing in Hip Arthroplasty

WSEAS TRANSACTIONS on SIGNAL PROCESSING, Volume 6, Print ISSN: 1790-5052, E-ISSN: 2224-3488, 2010

[MMM10c] *Anca Morar, Florica Moldoveanu, Alin Moldoveanu, **Victor Asavei***

Computer Assisted Analysis of 2D/3D Medical Images

Annals of DAAAM for 2010 & Proceedings of the 21st International DAAAM Symposium, ISBN 978-3-901509-73-5, ISSN 1726-9679, Editor B. Katalinic, Published by DAAAM International, Vienna, Austria 2010  
**(Proceedings ISI)**

[EMM10] *Egner, A. I.; Moldoveanu, F.; Moldoveanu, A. D. B.; **Asavei, V.**; Morar, A. - A.; & Boiangiu, C. A.*

Testing the interoperability of HL7-based applications using TTCN-3

Annals of DAAAM for 2010 & Proceedings of the 21st International DAAAM Symposium, ISBN 978-3-901509-73-5, ISSN 1726-9679, Editor B. Katalinic, Published by DAAAM International, Vienna, Austria 2010  
**(Proceedings ISI)**

[MAM10] *Moldoveanu, A. D. B.; **Asavei, V.**; Moldoveanu, F; Morar, A. - A. & Egner, A. I. .; Boiangiu, C. A.*

Turning Nearshoring into a Success – Managing Technical Background Differences  
Annals of DAAAM for 2010 & Proceedings of the 21st International DAAAM Symposium, ISBN 978-3-901509-73-5, ISSN 1726-9679, Editor B. Katalinic, Published by DAAAM International, Vienna, Austria 2010  
**(Proceedings ISI)**

[BPM10] *Boiangiu, C. A.; Petrescu, S.B. ; Moldoveanu, A. D. B.; .; **Asavei, V.**; Bucur, I.*

Systematic Printing Space Recognition  
Annals of DAAAM for 2010 & Proceedings of the 21st International DAAAM Symposium, Editor B. Katalinic, Published by DAAAM International, Vienna, Austria 2010 **(Proceedings ISI)**

[APM11] ***Victor Asavei**, Alexandru-Lucian Petrescu, Florica Moldoveanu*

A Fractal Approach to Terrain Generation and Rendering  
CSCS 18 - The 18th International Conference On Control Systems And Computer Science, 24-27 May 2011, Bucharest

[EMM11] *Alexandru Egner, Florica Moldoveanu, Alin Moldoveanu, Victor Asavei, Anca Morar*

Automated generation of TTCN-3 type set used for testing of healthcare applications  
CSCS 18 - The 18th International Conference On Control Systems And Computer Science, 24-27 May 2011, Bucharest

[MMM11] *Anca Morar, Florica Moldoveanu, Alin Moldoveanu, Victor Asavei, Alexandru Egner*

Computer Assisted Insertion of Prostheses Based on Medical Images  
CSCS 18 - The 18th International Conference On Control Systems And Computer Science, 24-27 May 2011, Bucharest

### **8.3 ALTE ACTIVITĂȚI ȘTIINȚIFICE REALIZATE ÎN PERIOADA DESFĂȘURĂRII DOCTORATULUI**

Proiect de cercetare PNCDII- Parteneriate: “**Sisteme incorporate tip neuroproteză pentru recuperarea persoanelor cu handicap neuromotor**” (SINPHA), membru în echipa de cercetare (2007-2010)

Proiect de cercetare EUREKA: "**Testarea fiabilității sistemelor medicale**"  
**ReTeMes**, membru în echipa de cercetare (2008-2010)

Proiect de cercetare Leonardo da Vinci : "**Strategic Planning for Sustainable Clustering of Collaborative SMEs**" (**SMECluster**), membru în echipa de cercetare (2008-2010)

Proiect de cercetare Leonardo da Vinci : "**Establishment of Sustainable Collaborative SME Networks**" (**SMENet**), membru în echipa de cercetare (2008-2010)

Proiect de cercetare PNCDII- Parteneriate: "**Sistem informatic avansat, bazat pe imagistică medicală, pentru producerea implanturilor personalizate dedicate artroplastiei de șold**" (**SABIMAS**), membru în echipa de cercetare (2008-2011)

Proiect de cercetare EuroStars Inovare: "**Risk Detection in Laboratory Information Systems**" (**RELIS**), membru în echipa de cercetare (2010-prezent)

Proiect de cercetare EuroStars Inovare: "**Enterprise Unified Guideline Engine**" (**EUGEN**), membru în echipa de cercetare (2010-prezent)

Proiect de cercetare EuroStars Inovare: "**Visualization of Patient Data for easy management of care processes**" (**VISUAL-D**), membru în echipa de cercetare (2011-prezent)

Proiect de cercetare EUREKA: "**Identificarea riscurilor medicale operaționale și detecția fraudelor**" (**MORIS FD**), membru în echipa de cercetare (2011-prezent)

Organizare Workshop „Graphics and Virtual Reality Summer Workshop”  
2008,2009,2010,2011 în cadrul Facultății de Automatică și Calculatoare

## 9 BIBLIOGRAFIE

- [MIL10] John L. Miller, John Crowcroft  
The near-term feasibility of P2P MMOG's  
NetGames '10 Proceedings of the 9th Annual Workshop on Network and Systems Support for Games
- [LAK09] Geetika T. Lakshmanan, Yuri G. Rabinovich, Opher Etzion  
A stratified approach for supporting high throughput event processing applications  
Proceedings of the Third ACM International Conference on Distributed Event-Based Systems, 2009
- [MAR09] Alma Martinez G., Hector Orozco A., Felix Ramos C., Mario Siller  
A Peer-To-Peer Architecture for Real-Time Distributed Visualization of 3DCollaborative Virtual Environments  
2009 13th IEEE/ACM International Symposium on Distributed Simulation and Real Time Applications
- [CHE08] Che S. , Others  
A Performance Study of General Purpose Applications on Graphics Processors  
Journal of Parallel and Distributed Computing, Volume 68, Issue 10, Pages 1370-1380, 2008
- [NVI09a] NVidia Corporation  
NVIDIA OptiX RayTracing Engine,  
SIGGRAPH 2009, NEW ORLEANS—Aug. 4, 2009
- [CAR11] John Carmack,  
GPU Race, Intel Graphics, Ray Tracing, Voxels and more!,  
[HTTP://WWW.PCPEER.COM/REVIEWS/EDITORIAL/JOHN-CARMACK-INTERVIEW-GPU-RACE-INTEL-GRAPHICS-RAY-TRACING-VOXELS-AND-MORE](http://www.pcp.com/reviews/editorial/john-carmack-interview-gpu-race-intel-graphics-ray-tracing-voxels-and-more)
- [BRE08] Breitbart J.  
Case studies on GPU usage and data structure design  
*Dept. of Computer Science and Electrical Engineering, Universitat Kassel, 2008*
- [FOL05] Foley T. and Sugerman J.  
Kd-tree acceleration structures for a GPU Raytracer  
*ACM SIGGRAPH/EUROGRAPHICS conference on Graphics hardware, ACM Press, Pages 15-22*
- [CHE06] Alvin Chen and Richard R. Muntz  
Peer Clustering: A Hybrid Approach to Distributed Virtual Environments

*The 5th Workshop on Network & System Support for Games 2006 — NETGAMES 2006*

[LAV02] Antonia Lucinelma Pessoa Albuquerque, Luiz Velho  
Togetherness through virtual worlds: How real can be that presence?  
În Fifth Annual International Workshop on Presence 2002, pages 435–447, 2002

[GAR03] Maia Garau  
The Impact of Avatar Fidelity on Social Interaction în Virtual Environments  
2003

[PET97a] S. Petrescu, Zoea Racoviță, G. Hera, F. Moldoveanu, Alin Moldoveanu  
Algorithm for logic operation on connected closed contours CSCS 11,  
Proceedings of the 11th Int. Conference on Control Systems and Computer Science, vol. II, pg.  
260-264, 1997

[WAL04] Wald, Ingo  
Real Time Ray-Tracing and Interactive Global Illumination  
2004

[MOL96] Moldoveanu, Florica & Alții  
Grafica pe calculator  
Ed. Teora, București, 1996

[BLINN] Blinn, James  
Models of light reflection for computer synthesized pictures  
SIGGRAPH '77 Proceedings of the 4th annual conference on Computer graphics and interactive  
techniques

[PHONG] Phong, B.T.  
Illumination for computer generated pictures  
Communications of ACM 18

[QUA04] Quake3 RayTraced  
[HTTP://WWW.Q3RT.DE](http://www.Q3RT.DE)

[QUA06] Quake4 RayTraced  
[HTTP://WWW.Q4RT.DE](http://www.Q4RT.DE)

[QWR08] Quake Wars RayTraced  
[HTTP://WWW.QWRT.DE](http://www.QWRT.DE)

[WOL10] Wolfenstein RayTraced  
[HTTP://WWW.WOLFRT.DE](http://www.WOLFRT.DE)



[GAI11] Gaikai  
[HTTP://WWW.GAIKAI.COM](http://www.gaikai.com)

[ONL10] OnLive  
[HTTP://WWW.ONLIVE.COM](http://www.onlive.com)

[EVE11] Eve OnLine  
[HTTP://WWW.EVEONLINE.COM](http://www.eveonline.com)

[WOW11] World of Warcraft  
[HTTP://WWW.WOW-EUROPE.COM](http://www.wow-europe.com)

[CUDA11] CUDA Programming Guide  
[HTTP://WWW.NVIDIA.COM/OBJECT/CUDA\\_DEVELOP.HTML](http://www.nvidia.com/object/cuda_develop.html)

[FEK06] Fernando Randima, Kilgard Mark  
The Cg Tutorial-The Definitive Guide to Programmable Real-Time Graphics  
Addison-Wesley, 2006

[FAT08] Kayvon Fatahalian and Mike Houston  
A Closer look at GPUs  
Communication of the ACM, vol 51, no 10, 2008

[WAO82] M. W. Walker, D. Orin  
Efficient dynamic computer simulation of robotic mechanisms  
ASME Journal of Dynamic Systems, Measurement, and Control, 1982

[FEA87] R. Featherstone  
Robot Dynamics Algorithms  
Kluwer Academic Publishers, 1987

[BMT96] Ronan Boulic, Ramon Mas, Daniel Thalmann  
A robust approach for the control of the center of mass with inverse kinetics  
Computers & Graphics 20(5): 693-701,1996

[PER02] K Perlin  
Generation based on motion blending  
Proceedings of ACM SIGGRAPH Symposium on Computer Animation, 2002]

[WIH97] D J Wiley, J K Hahn  
Interpolation synthesis for articulated figure motion

Proceedings of Virtual Reality Annual International Symposium, 1997

[GUR96] Shang Guo, James Robergé

**A high-level control mechanism for human locomotion based on parametric frame space interpolation**

Proceedings of the 6th EuroGraphics Workshop on Animation and Simulation, 1996

[RCB98] Charles Rose, Michael F. Cohen, Bobby Bodenheimer

**Verbs and Adverbs: Multidimensional Motion Interpolation**

IEEE COMPUTER GRAPHICS AND APPLICATIONS, Volume 18 Issue 5, 1998

[DEN04] Philippe Decaudin, Fabrice Neyret

**Rendering Forest Scenes in Real-Time**

RENDERING TECHNIQUES '04 ,EUROGRAPHICS SYMPOSIUM ON RENDERING, 2004

[CDN04] Florent Cohen, Philippe Decaudin, Fabrice Neyret

**GPU-Based Lighting and Shadowing of Complex Natural Scenes**

SIGGRAPH Conf., 2004

[CHA06] Chris Chambers, Wu-chang Feng, Wu-chi Feng

**Towards public server MMOs**

Proceedings of 5th ACM SIGCOMM workshop on Network and system support for games, 2006

[KES05] J. Kesselman

**Server Architectures for Massively Multiplayer Online Games**

in Session TS-1084, Javaone conference, SUN, 2005

[WIE06] A. Wierzbicki

**Trust Enforcement in Peer-to-peer Massive Multi-player Online Games**

Proc. Grid computing, high-performance and Distributed Applications (GADA'06), Springer, LNCS 4276 (2006), 1163-1180

[GUP09] Nitin Gupta and others

**Scalability for Virtual Worlds**

Proceedings ICDE '09 Proceedings of the 2009 IEEE International Conference on Data Engineering

[BEZ08] Carlos Eduardo B. Bezerra, Fábio R. Cecin, Cláudio F. R. Geyer

**A3: a Novel Interest Management Algorithm for Distributed Simulations of MMOGs**

12th 2008 IEEE/ACM International Symposium on Distributed Simulation and Real-Time Applications

- [GLI08] Frank Glinka, Alexander Ploss, Sergei Gorlatch, and Jens Müller-Iden  
**High-level development of multiserver online games**  
International Journal of Computer Games Technology - Networking for Computer Games, Volume  
2008, January 2008