

3D Ray Tracing

Mihai Gheoace 343C1

Cezar Crăciunoiu 343C1

Specificația cerințelor

Introducere

Proiectul 3D Ray Tracing are ca scop rularea algoritmului de ray tracing într-un spațiu tridimensional pe un număr limitat de obiecte. Deoarece acest algoritm este unul neperformant atunci când nu există suport hardware, am dorit să facem o comparație a rulării acestuia, atât pe CPU cât și pe GPU.

Pentru rasterizarea imaginilor se pune în evidență arhitectura SISD (la rulare s-a folosit OpenMP pentru a paraleliza ușor pe CPU, deoarece timpii de rulare erau foarte mari. Se poate alege totuși să se ruleze SISD în loc de SIMD) a procesorului în raport cu arhitectura SIMD a GPU-ului. La final se va face o comparație a performanței de rulare între cele două versiuni.

Descriere generală

Pentru rularea proiectului se utilizează un procesor Intel 7700k cu o placă grafică Nvidia GTX 1060. La nivel software am utilizat Visual Studio Enterprise 2019 și Framework-ul de laborator oferit de echipa de SPG. Am folosit versiunea 440 de OpenGL pentru rularea shader-urilor.

Am ales acest proiect pentru a sublinia rolul crucial al plăcii grafice în execuția unor proiecte ce folosesc intensiv resursele calculatorului într-o paradigmă puternic paralelizabilă.

Sistemul propus

Algoritmul poate fi asimilat în alte proiecte care sunt construite în Framework-ul de laborator SPG. Cerințele de sistem ale proiectului sunt un procesor ce suporta versiunea 440 de OpenGL cât și Framework-ul de laborator care sa fie configurat corespunzător rulării. De asemenea, pentru rularea Framework-ului cu Visual Studio este nevoie ca sistemul de operare să fie Windows.

La general, proiectul va afișa trei scene de complexitate diferită, ce vor fi procesate atât pe GPU cât și pe CPU. La comanda utilizatorului se va face trecere de pe CPU pe GPU și invers. Totodată, acesta va putea să treacă prin scene folosind tastele numerice.

Caz de utilizare

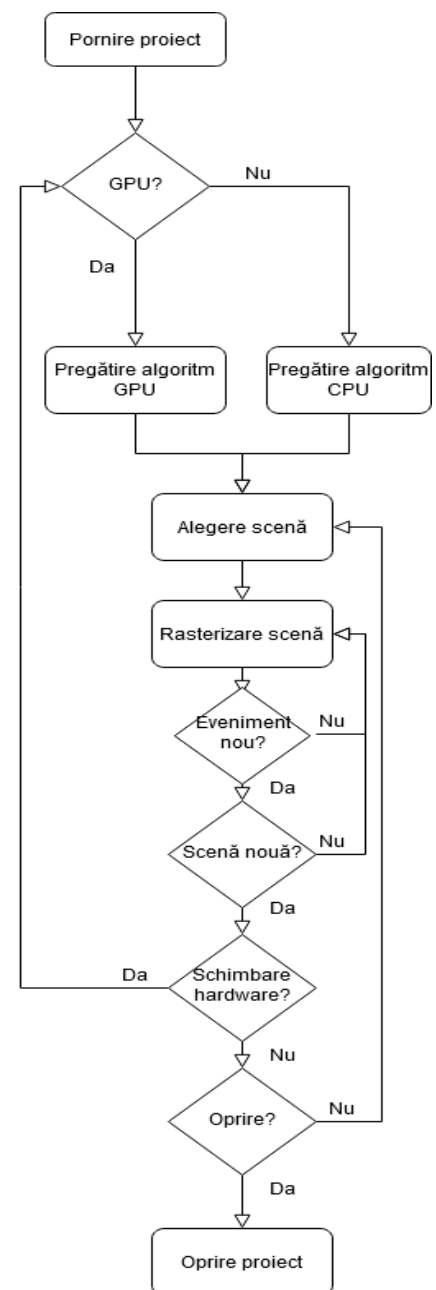
Precondiții:

- Visual Studio 2019 (minim) instalat pe Windows (> 7)
- Framework de laborator bine configurat
- Proiect construit

Postcondiții:

- Se observă în timp real diferența dintre CPU și GPU

În organigrama alăturată se observă fluxul proiectului. Utilizatorul pornește proiectul și alege dacă să ruleze pe GPU sau CPU. Programul realizează în spate pregătirea algoritmului și apoi utilizatorul trebuie să aleagă una din cele trei scene pentru a le reda. Scena este afișată în timp real, până când utilizatorul apasă o tastă numerică pentru a alege o scenă nouă sau una din tastele '-' și '+' pentru a alege CPU, respectiv GPU. Dacă se apasă pe 'X' programul se oprește.



Arhitectura

Prezentare generala sumară

Proiectul este făcut să meargă pe orice calculator ce poate rula framework-ul de laborator SPG, implicit Visual studio 2019 (cu OpenGL 440). De asemenea, o placă video este puternic recomandată pentru rularea decentă a proiectului. Nu sunt necesare alte elemente adiționale pentru execuția proiectului.

Decompoziția în subsisteme

După cum s-a precizat, proiectul este structurat în două versiuni, una ce rulează pe CPU și una ce rulează pe GPU. Astfel, și implementările sunt diferite. Cele două încep comun, cu inițializarea programului folosind framework-ul de laborator. Tot în mod comun se realizează cronometrarea și contorizarea cadrelor pe secundă, respectiv interacționarea cu utilizatorul. Sistemele se despart în partea de implementare propriu-zisă, o variantă fiind scrisă în limbajul GLSL, iar alta direct în CPP. Putem vorbi astfel de un subsistem pentru GPU și unul pentru CPU.

Cazuri de utilizare limită

În anumite cazuri, cum ar fi atunci când camera este poziționată incorect în interiorul obiectelor, sistemul nu mai poate trasa corespunzător razele, deoarece acestea nu ajung niciodată în interior. Astfel, se ajunge la o rulare la capacitate maximă, chit că nu ar fi nevoie.

Detalii de implementare

Discuția în legătură cu implementarea va fi împărțită în două. Chit că ambele implementări fac același lucru, acestea se aseamănă doar conceptual. Atât limbajul, cât și abordarea sunt diferite. De asemenea, deoarece generarea pe procesor este foarte lentă, măsurarea pe acesta se va face în cadre pe oră, pentru a obține valori apropiate de cele de pe GPU (chit că unitatea de măsură este diferită).

Implementare GPU

Pentru această implementare, s-a trimis către rasterizare de pe CPU un singur quad. Acesta a fost amplasat folosind coordonate dispozitiv normalizate în fața camerei în Vertex Shader. Astfel încât, putem spune că, algoritmul Ray Tracing generează tot imagini și le aplică pe quad, dar, deoarece acest lucru se face foarte repede, se obține ideea de video. Tot în Vertex Shader se calculează 2 valori *eye_pos* și *ray_direction_from_vertex*, ce vor fi folosite pentru trasarea corectă a razelor.

Toate calculele sunt realizate în Fragment Shader. Pentru a se realiza obiectele, s-au folosit mai multe structuri, prin intermediul cărora s-au definit următoarele forme geometrice: paralelipiped/cub, sferă, cilindru, con. Obiectele nu urmează pipeline-ul clasic OpenGL, fiind generate direct în Fragment Shader. Acestea reprezintă zone de intersecție cu razele și absorb diferite lungimi de undă (culori) din lumină. Ce rămâne este culoarea obiectelor.

În rest, se urmează algoritmul clasic de trasare a razelor. Se calculează intersecțiile, se verifică dacă un obiect este în umbra altui obiect, se calculează lumina Phong de tipul ambientală/difuză/speculară, și, în final, se verifică dacă obiectul este transparent sau reflectorizant pentru a se modifica razele.

Mai jos sunt discutate funcțiile ce se ocupă cu intersecția obiectelor cu razele. Aceasta este partea cea mai intensivă a codului. Toate aceste funcții întorc cel mai apropiat punct de intersecție cu raza.

Intersecție Paralelipiped

Se folosește algoritmul Axis-aligned Bounding Boxes (AABB) pentru a genera ușor și rapid un paralelipiped (inclusiv Cub). Aceasta este, de departe, cea mai rapidă și ușoară metodă de intersecție a obiectelor.

Intersecție Sferă

Pentru a calcula intersecția cu sfera s-a utilizat ecuația de gradul 2 prin care s-au obținut 2 rădăcini. Dintre acestea se alege cea mai mare. Această intersecție este ceva mai înceată, dar, recuperează prin faptul ca este succintă.

Intersecție Cilindru

Pentru generarea cilindrului s-a urmat modelul sferei. Mai întâi se generează capacele cilindrelor sub formă de cercuri și apoi se folosește ecuația de gradul 2 în funcție de X și Y, pentru a se genera un cilindru infinit, care apoi este crestat la capete în funcție de dimensiuni.

Intersecție Con

Se urmează modelul cilindrului, dar ecuația de gradul 2 este ridicată la pătrat pentru a genera planurile în formă de X. Apoi se decupează forma până la dimensiunea dorită și se pune capacul conului. De asemenea se elimină jumătatea mai mare decât centrul din figură, pentru a se înjumătăți.

Implementare CPU

Codul de pe CPU este împărțit în mai multe fișiere, întrucât are o structură OOP, plecând de la propunerile din cartea Ray Tracing in One Weekend de Peter Shirley. Funcțiile de calculare a unui hit între o rază și un obiect sunt similare cu abordarea de pe GPU(deci intersecțiile), însă observăm câteva diferențe de abordare:

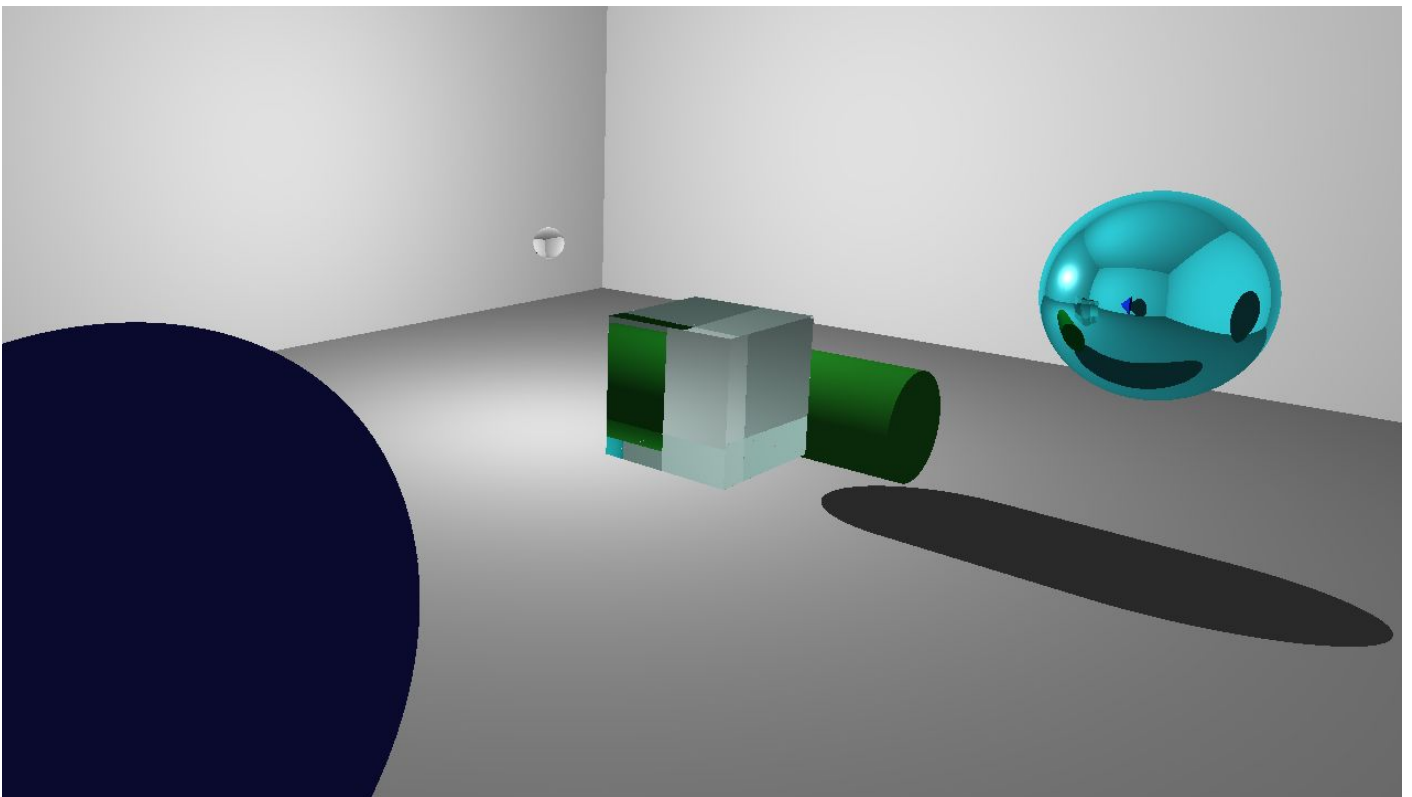
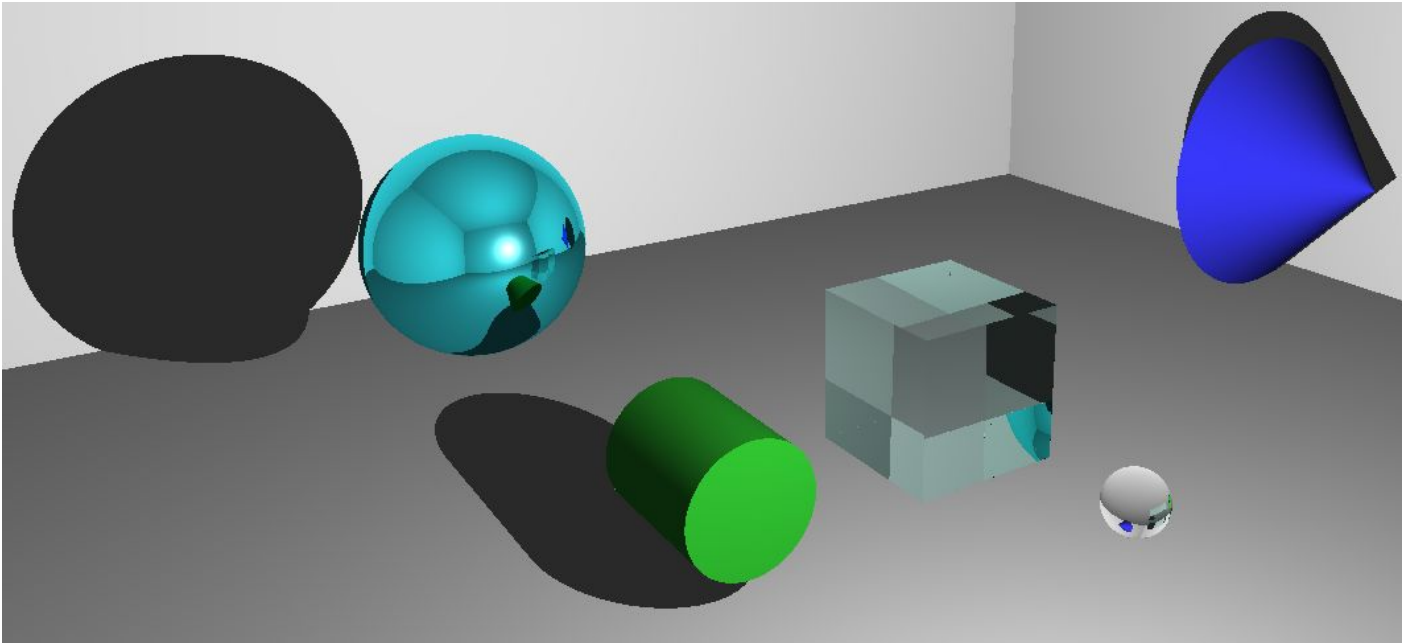
- toate obiectele sunt parcurse în aceeași buclă, în fiecare pas calculându-se normala, care este orientată invers față de sensul razei incidente (pe GPU este întotdeauna spre exteriorul obiectului).
- normala este calculată în același timp în care se descoperă cel mai apropiat obiect de originea razei.
- se realizează antialiasing prin realizarea de multiple calcule per pixel prin for-ul cel mai intern din funcția Ray.
- codul este paralelizat pe CPU prin openMP(versiunea 2.0 de openMP pentru Visual Studio 2019) - se obține speedup de ~3.5 pentru 4 core-uri (8 logice) față de versiunea serială ceea ce folosește la maxim procesorul. Acest pas este opțional, dar, în cele mai multe cazuri timpii de rulare sunt foarte mari, deci se recomandă paralelizarea.
- se declară adâncimea de vizualizare a camerei explicit.

Screenshots

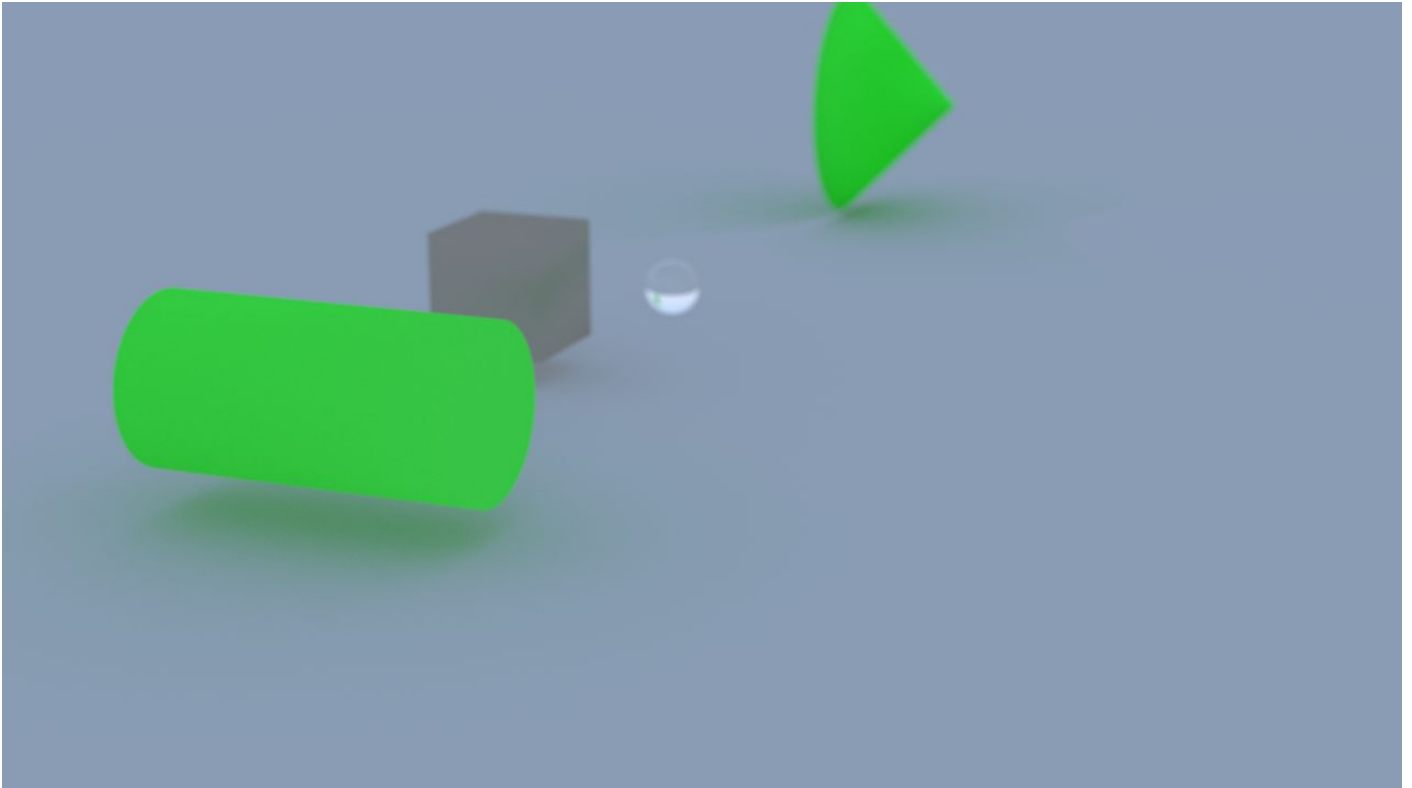
Scena 1

Scena 1 este folosită pentru a exemplifica obiectele și proprietățile acestora într-o lume ce folosește RayTracing.

GPU



CPU

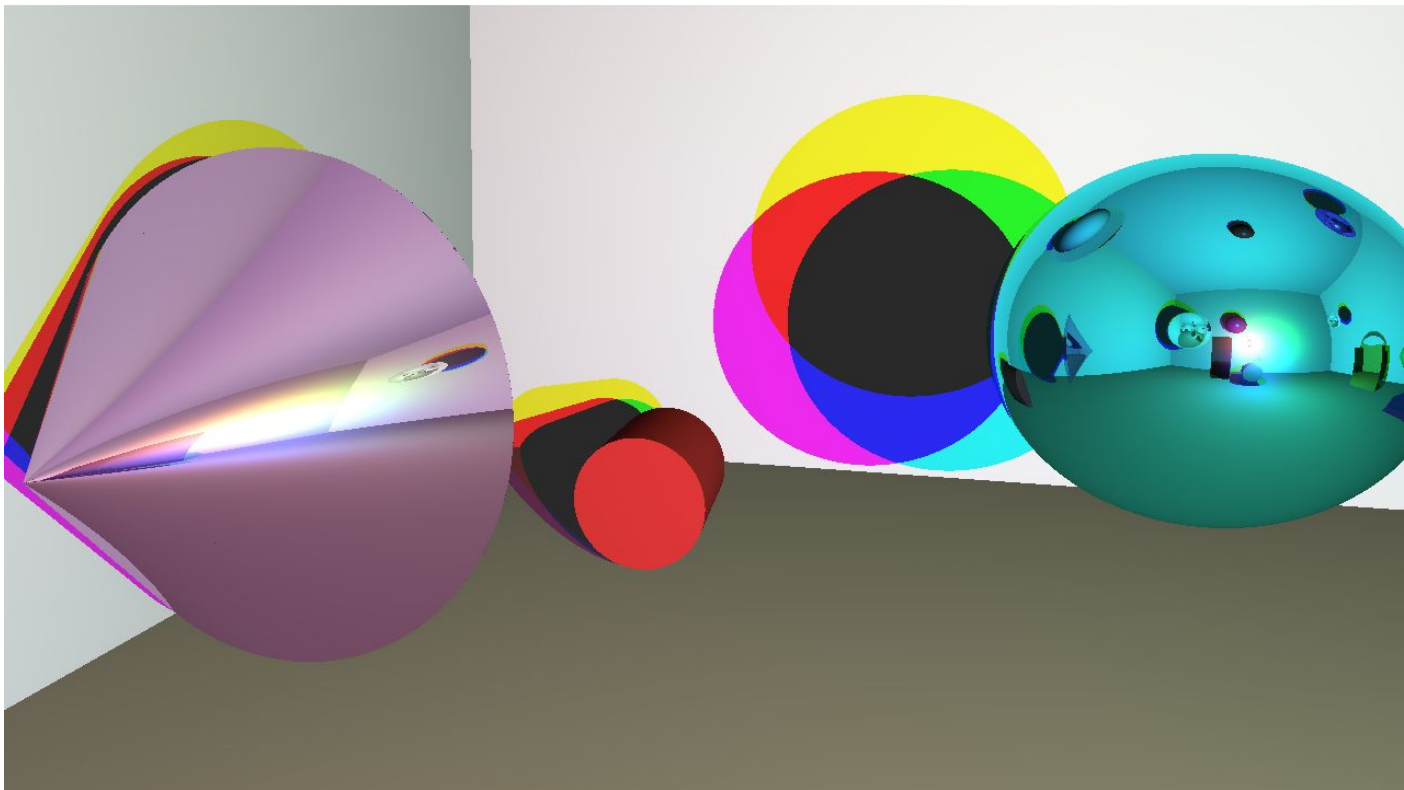
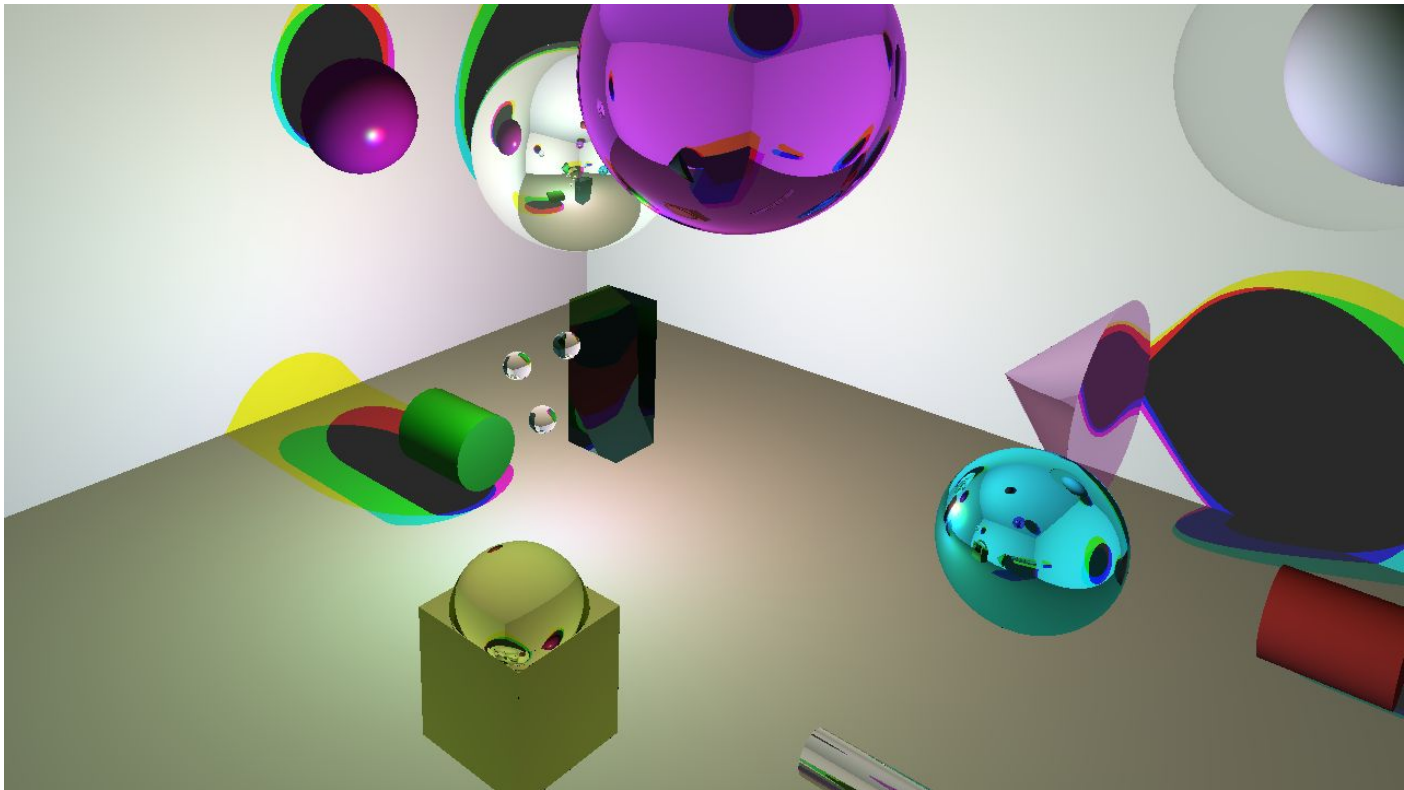


Scena 2

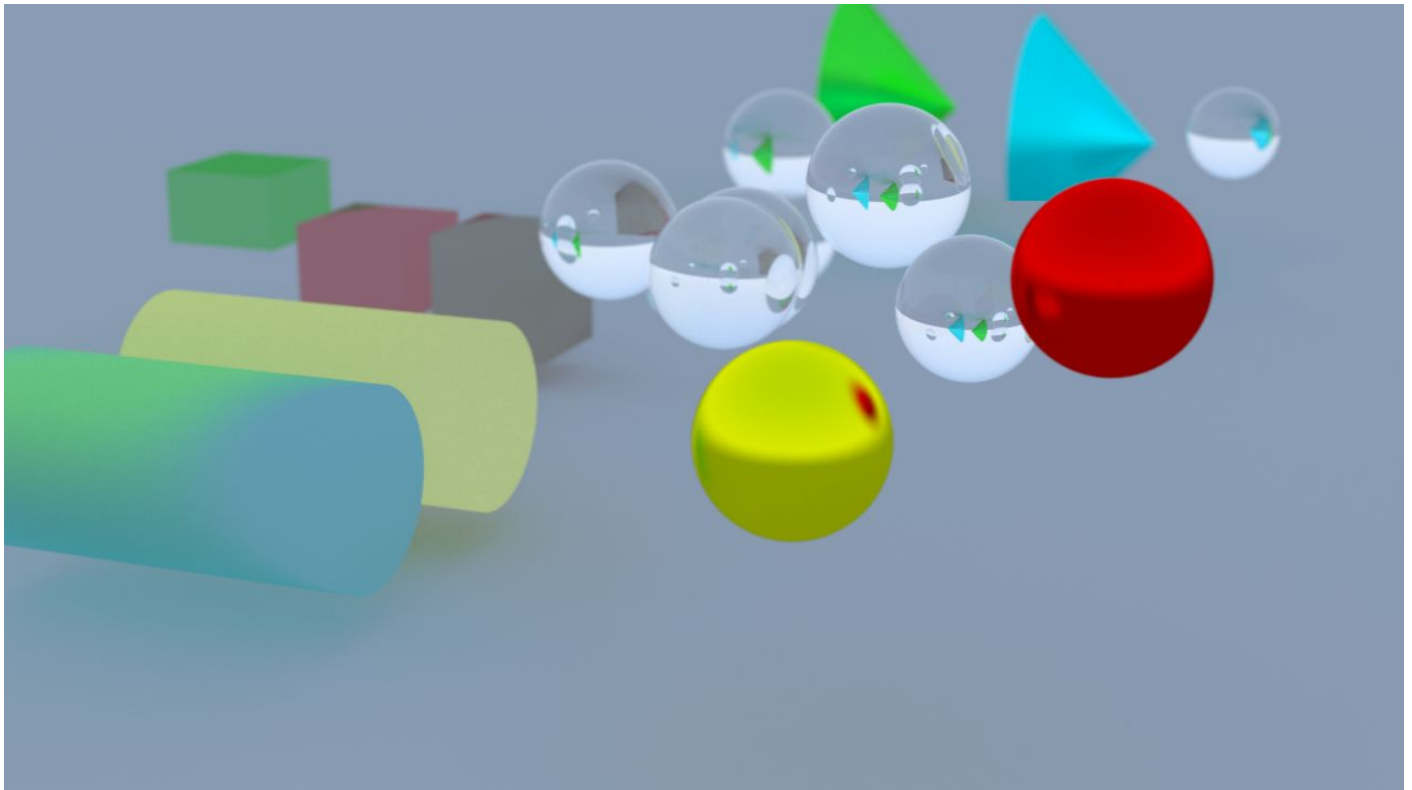
Această scenă încarcă majoritatea combinațiilor de obiecte și 3 surse de lumină, pentru a genera o lume standard ca dificultate de randare.

GPU

Sursele de lumină sunt de tip R, G și B. Astfel că, deoarece umbra a fost calculată ca opusul luminii, se obțin umbre în cercul de culori CMYK. Deși nu foarte corect, acest lucru a fost păstrat pentru aspectul vizual plăcut și rezultatele interesante.



CPU

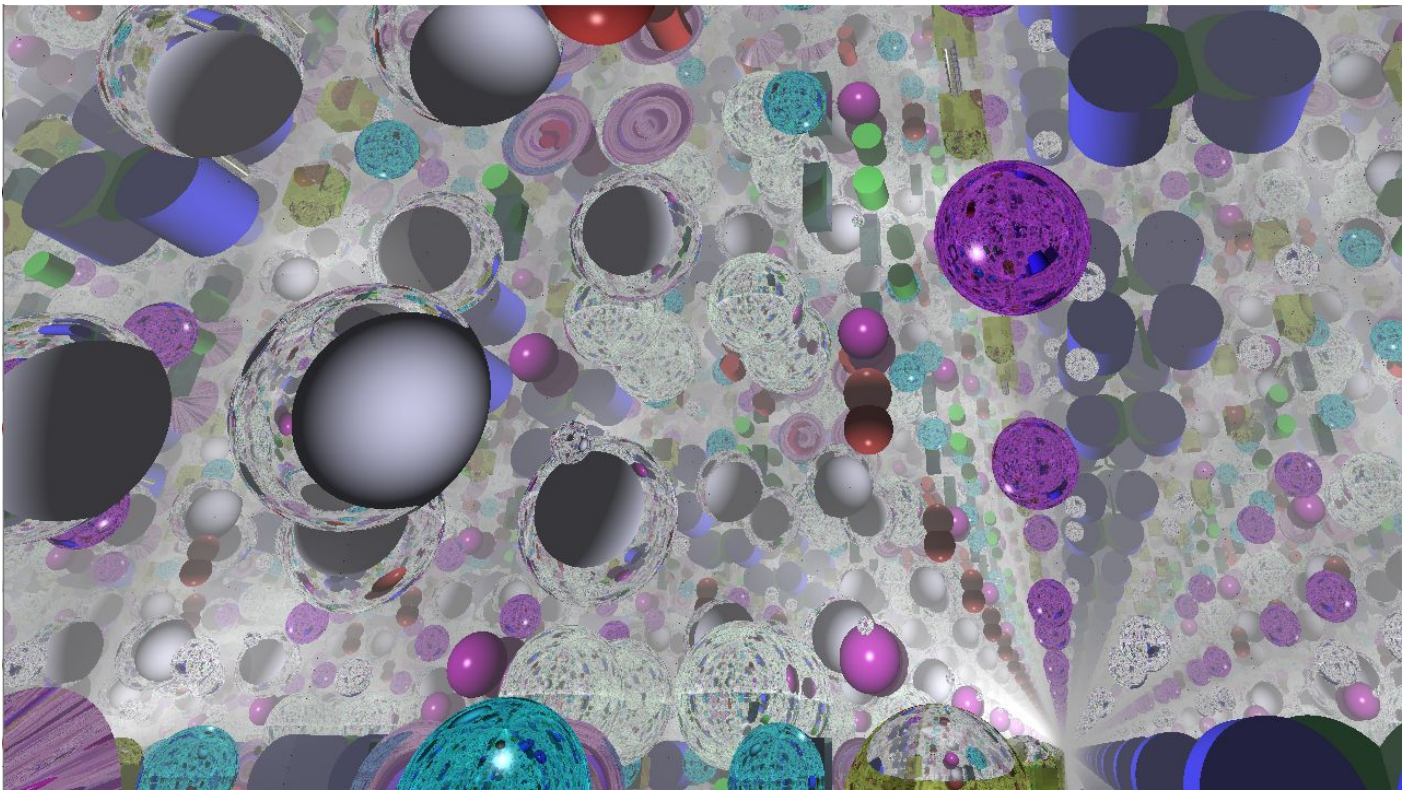
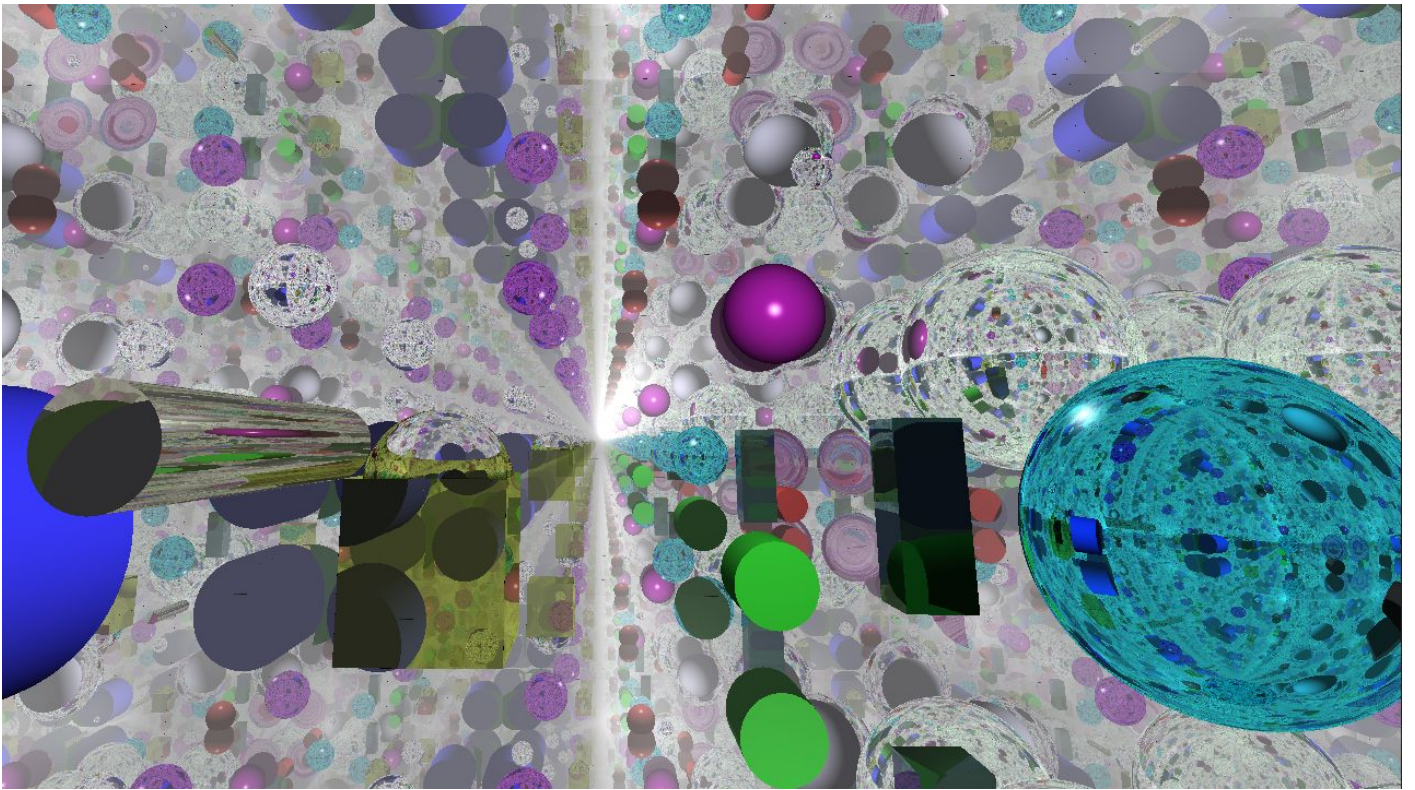


Pe CPU este pus accentul pe numărul mare de obiecte cu reflexii sub ele realiste și ușor blur în fundal pentru efectul de depth-of-field.

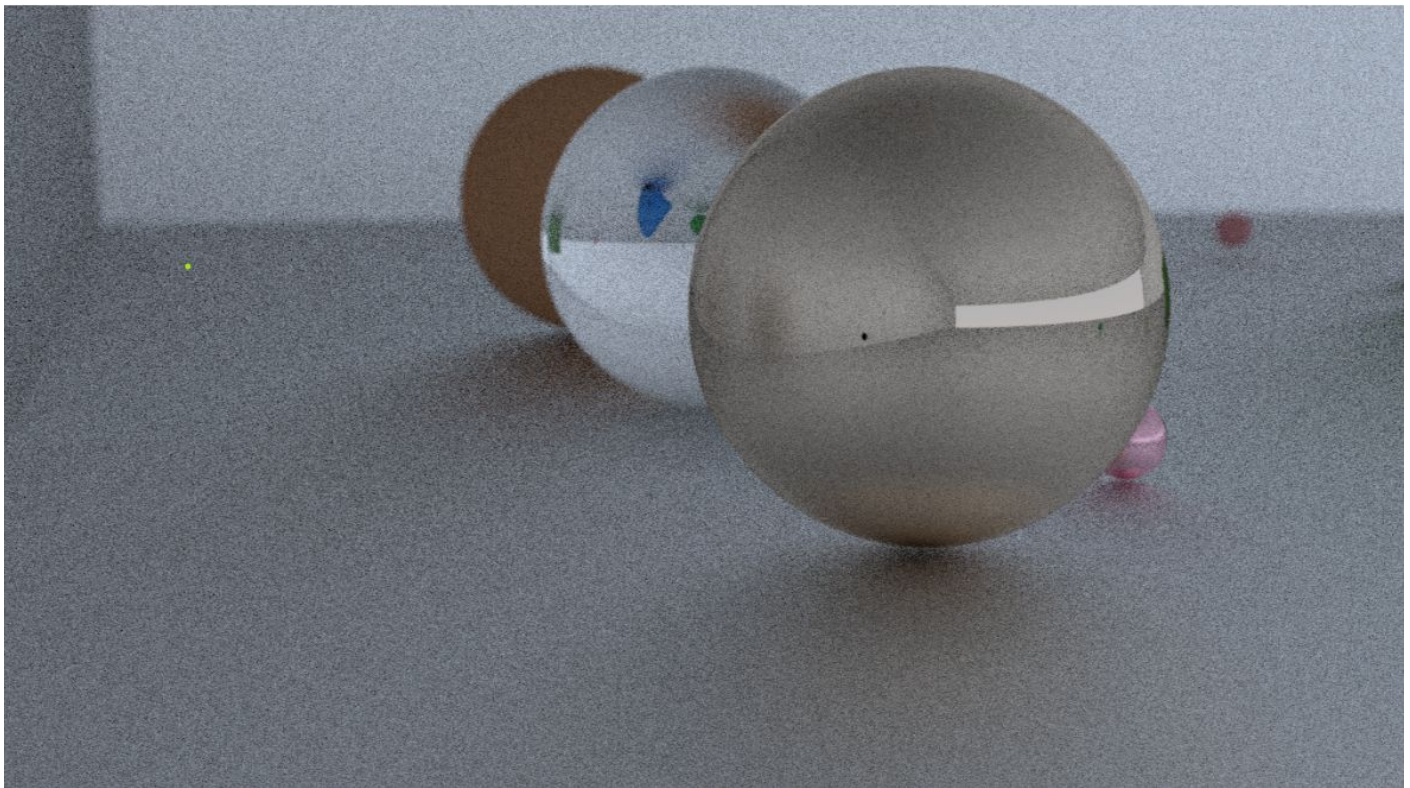
Scena 3

La fel de multe obiecte ca în a doua scenă, dar, mediul înconjurător reflectă lumina aproape perfect. Se obține astfel un efect interesant și foarte consumator de resurse.

GPU



CPU



Testare

Pentru a testa rularea programului am verificat rularea celor 3 scene și am comparat încărcarea CPU/GPU-ului, respectiv FPH/FPS-ul în timpul rulării. Pe GPU s-au folosit 40 de bounce-uri maxime.

Se obțin următoarele măsurători instantanee pentru cele 3 scene.

Idle GPU

GPU

NVIDIA GeForce GTX 1060 6GB

3D

0%

Copy

0%

Video Encode

0%

Video Decode

0%

Dedicated GPU memory usage

6.0 GB

Shared GPU memory usage

8.0 GB

Utilization

0%

Dedicated GPU memory

0.9/6.0 GB

GPU Memory

1.0/14.0 GB

GPU Temperature

55 °C

Driver version:

27.21.14.6079

Driver date:

2020-12-03

DirectX version:

12 (FL 12.1)

Physical location:

PCI bus 1, device 0, function 0

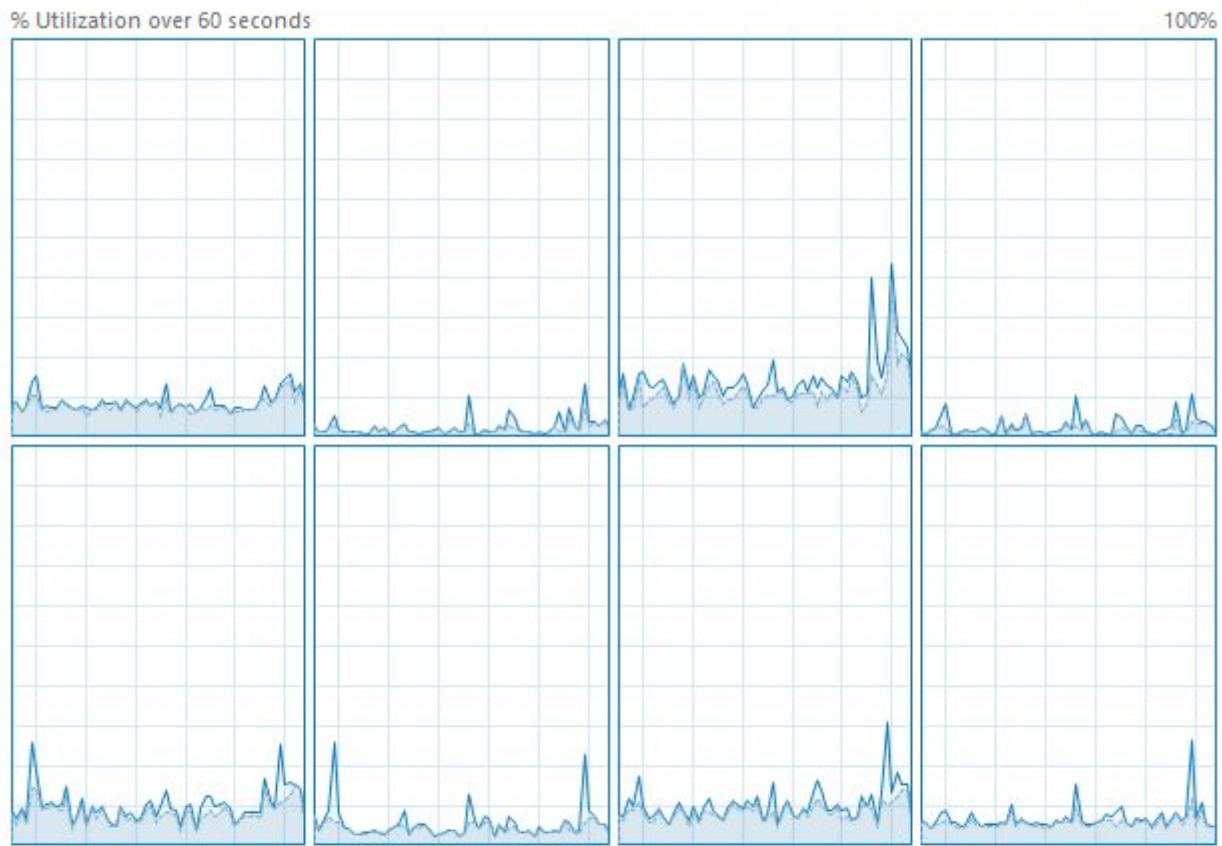
Hardware reserved memory:

92.0 MB

Idle CPU

CPU

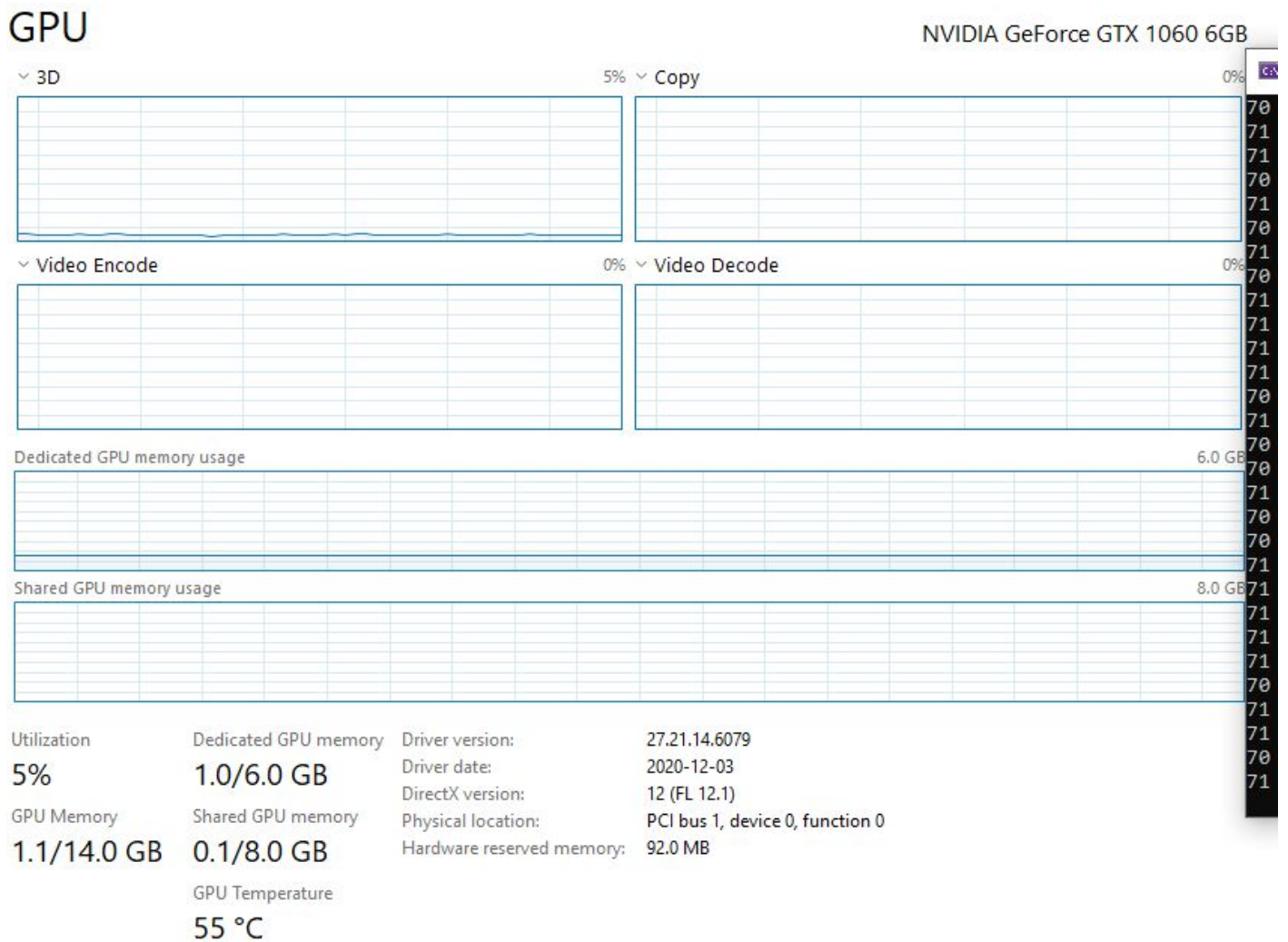
Intel(R) Core(TM) i7-7700 CPU @ 3.60GHz



Utilization	Speed		Base speed:	3.60 GHz
6%	4.09 GHz		Sockets:	1
			Cores:	4
Processes	Threads	Handles	Logical processors:	8
231	2772	98490	Virtualization:	Enabled
Up time			L1 cache:	256 KB
0:00:56:52			L2 cache:	1.0 MB
			L3 cache:	8.0 MB

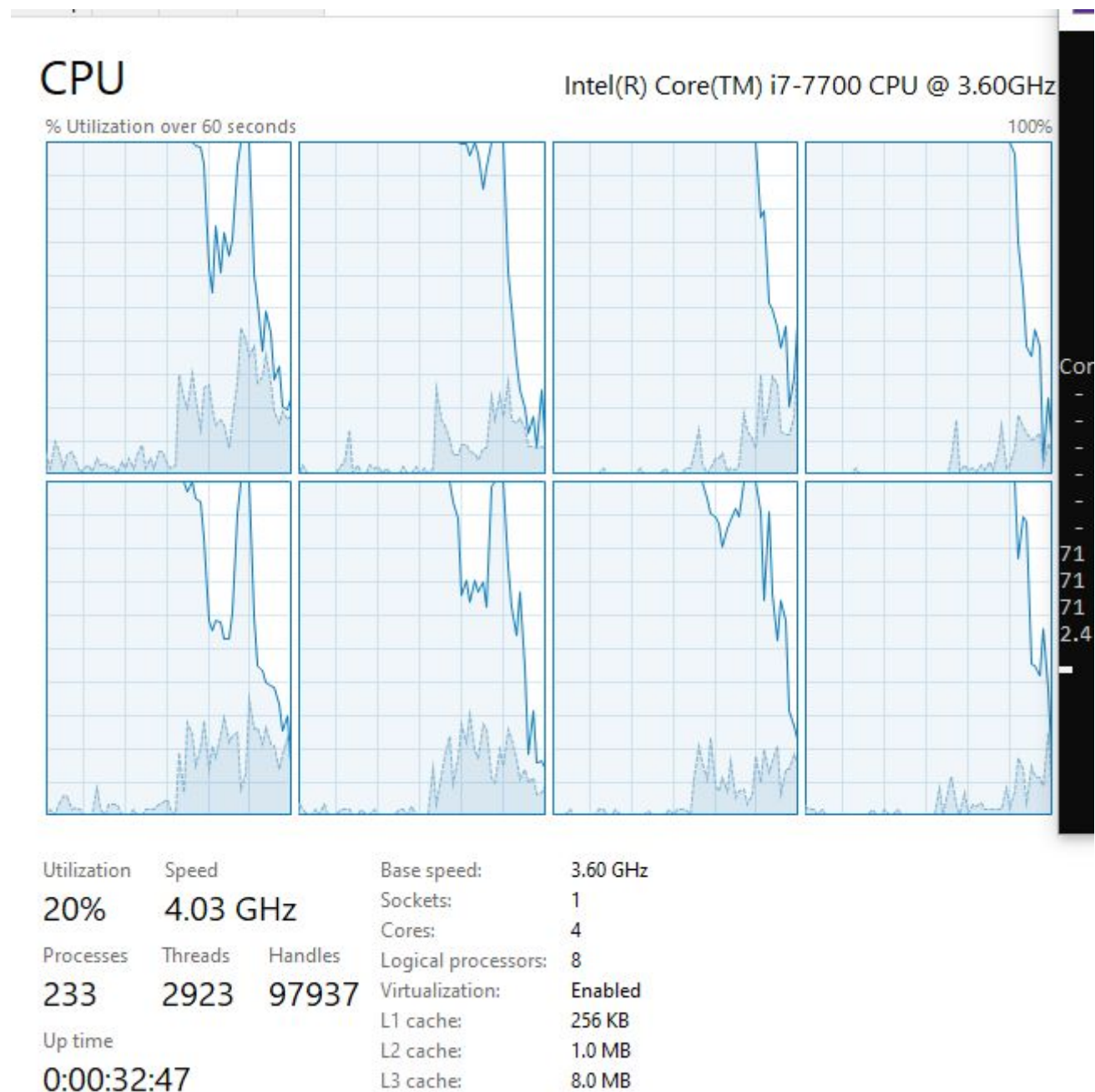
Scena 1

GPU



Se observă o încărcare de 5% și FPS-uri maxime (monitor de 71 hz).

CPU



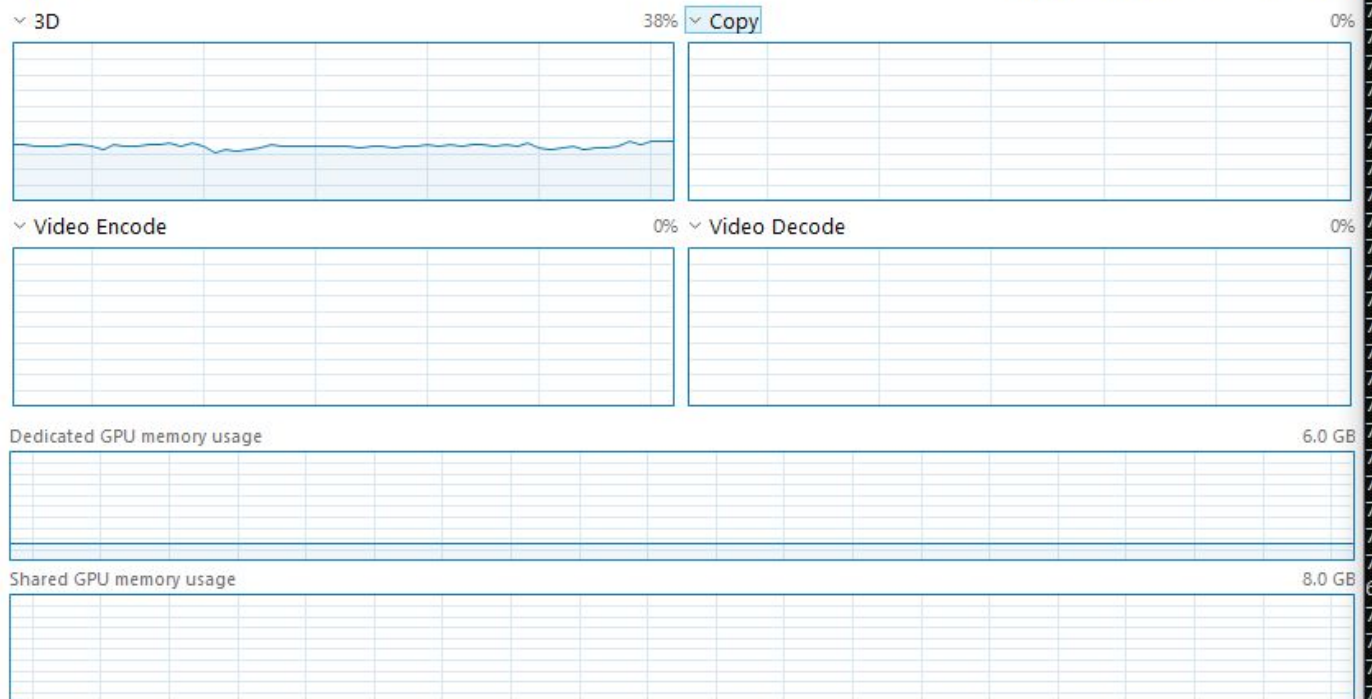
Se observă o rată de generare de 2.4 FPH. Sample rate-ul a fost de 600 și o adâncime maximă de 100.

Scena 2

GPU

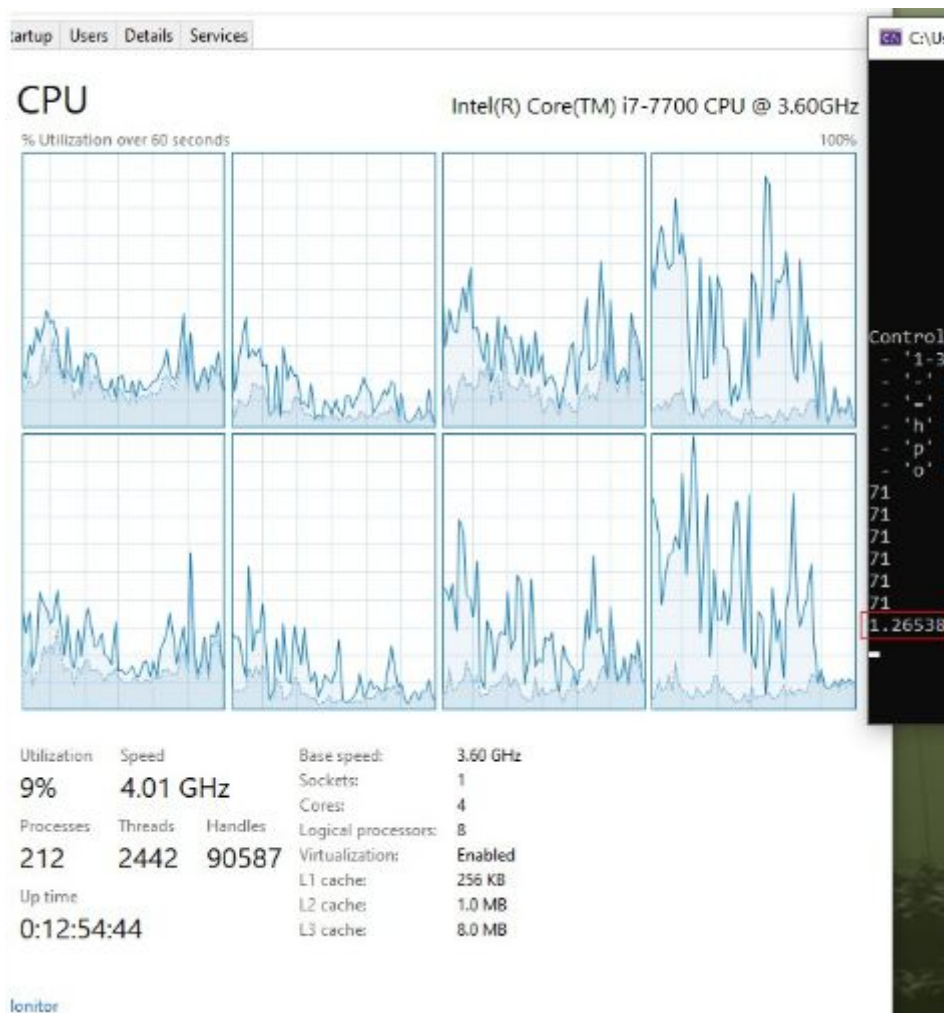
GPU

NVIDIA GeForce GTX 1060 6GB



Se observă o încărcare de ~40% și FPS-uri în continuare maxime.

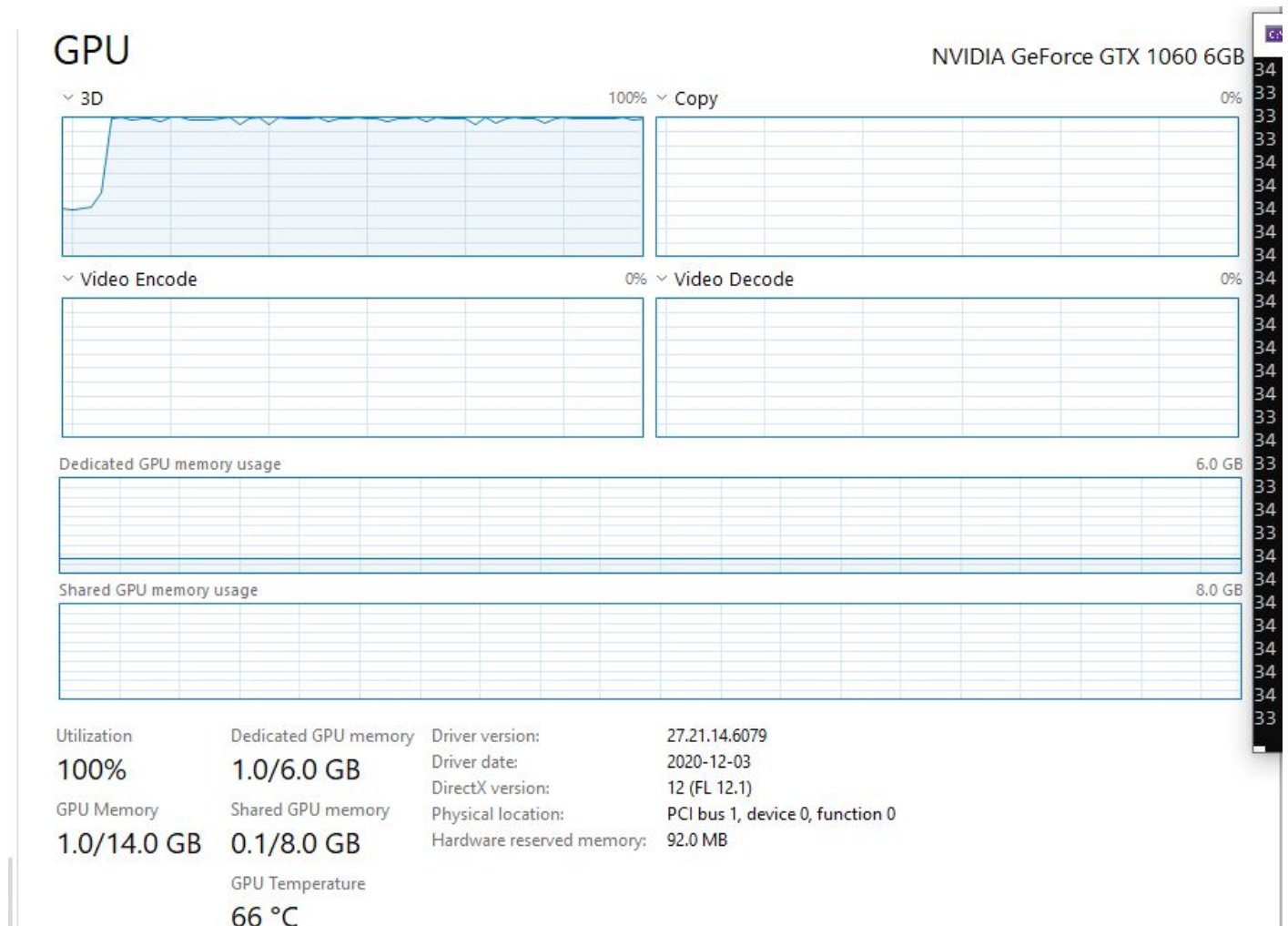
CPU



Încărcarea procesoarelor nu este mare deoarece captura a fost realizată ulterior. Totuși, se observă numărul mic de FPH. Sample rate-ul folosit a fost de 400 și distanța maximă de 100.

Scena 3

GPU



Se observă o încărcare maximă și FPS-uri precare ~34 (s-a rulat cu 40 de reflexii maxim a razelor).

CPU

CPU

Intel(R) Core(TM) i7-7700 CPU @ 3.60GHz

% Utilization over 60 seconds

100%



Utilization	Speed		Base speed:	3.60 GHz
10%	4.02 GHz		Sockets:	1
			Cores:	4
Processes	Threads	Handles	Logical processors:	8
251	3099	108387	Virtualization:	Enabled
Up time			L1 cache:	256 KB
0:03:45:18			L2 cache:	1.0 MB
			L3 cache:	8.0 MB

S-a rulat cu 15 sample per pixel si 10 pentru max_depth. Aceste valori s-au ales deoarece, pentru configurația clasică de 400 cu 100 s-ar fi obținut $8 / 1333.33$ FPH.

Bibliografie

1. Ray Tracing in One Weekend , Peter Shirley, *Version 3.2.3, 2020-12-07*,
<https://github.com/RayTracing/raytracing.github.io>