

# APLICAȚII PENTRU PROCESAREA DE IMAGINI PE PLATFORMA CUDA – STUDIU DE CAZ

**Ramona Din**

din.ramona@yahoo.com

**Andrei Grumăzescu**

andrei.grum@yahoo.com

**Daniela Saru**

saru@aii.pub.ro

**Ștefan Mocanu**

smocanu@rdslink.ro

**Radu Dobrescu**

radud@isis.pub.ro

Universitatea “Politehnica” București

**Rezumat:** Procesoarele grafice (GPUs – *Graphical Processing Units*) au cunoscut în ultimii ani o evoluție spectaculoasă, apropiindu-se din ce în ce mai mult de statutul de co-procesor al unității centrale de prelucrare. Pe lângă funcțiile specializate de prelucrare/redare a imaginilor, plăcile video oferă în prezent facilități performante pentru realizarea calculului matematic complex și rularea unor aplicații ce necesită putere ridicată de procesare. Exploatarea platformelor asociate plăcilor video a devenit accesibilă tuturor programatorilor odată cu lansarea arhitecturilor dedicate GPGPU. Această lucrare își propune evidențierea performanțelor pe care le oferă una dintre cele mai competitive arhitecturi ale momentului, CUDA, pentru elaborarea unor algoritmi optimizați de prelucrare a imaginilor. Studiul de caz prezintă rezultatele obținute în urma implementării unui algoritm de calcul paralel al histogramei unei imagini pe 256 niveluri de gri și o analiză comparativă a performanțelor obținute în urma prelucrării pe CPU, respectiv pe GPU.

**Cuvinte cheie:** procesor grafic (GPU), CUDA, procesare paralelă, GPGPU, histogramă.

**Abstract:** Graphic processors (GPUs), originally designed exclusively for game industry and visual effects, have turned into complex multi-threaded architectures which are now closer than ever to gain the co-processor status. GPUs can be used to process data and to deal with computationally intensive applications really fast. Programmers' access to these platforms has been enhanced along with the official launch of the GPGPU dedicated architectures. This paper focuses on the study of the performances offered by one of the most competitive architecture, CUDA, used to elaborate optimized image processing algorithms. This case-study presents the results gathered running a parallel algorithm designed to compute a 256 gray level image histogram on CUDA platform and it points out the performance differences between the results collected when the algorithm was run on the CPU vs. on the GPU platform.

**Key words:** Graphical processing unit (GPU), CUDA, parallel computing, GPGPU.

## 1. Introducere

În contextul cerințelor pieței pentru o putere de calcul din ce în ce mai mare, dar și ca urmare a evoluției din punct de vedere tehnologic, prelucrarea de tip secvențial a informației a fost înlocuită cu algoritmi proiectați pentru rularea în paralel a mai multor fire de execuție, în special în cazul volumelor mari de date ce necesită prelucrări complexe. Dacă în urmă cu câțiva ani exploatarea puterii de calcul a unui sistem se rezuma la paralelizarea algoritmilor pentru nucleele procesoarelor unității centrale, în ultimul timp s-a urmărit proiectarea unor arhitecturi hardware care să permită procesare în paralel și în afara acestora. Astfel, a apărut conceptul *General-Purpose computation on Graphics Processing Units (GPGPU)* [4]. Spre deosebire de arhitectura procesorului central, procesorul grafic (GPU) a fost proiectat pentru execuția foarte rapidă a unui set redus de instrucțiuni, având la bază mai multe nuclee de procesare pe care pot rula în paralel diferite taskuri. În timp ce în cazul CPU se acordă o atenție deosebită accesului la memorie sau transferului de date, proiectarea GPU a pus accentul pe optimizarea operațiilor ce implică procesarea efectivă a datelor [1]. Mai mult, trebuie menționat faptul că memoria video GDDR3 (accesată și utilizată de GPU pentru stocarea informațiilor) este mult mai rapidă decât memoria unității centrale, DDR2 sau DDR3.

Odată cu dezvoltarea puterii de calcul din punct de vedere tehnologic, producătorii de pe piața componentelor hardware au înțeles că este necesară crearea unor platforme care să permită exploatarea puterii GPU. Soluția celor de la NVIDIA pentru GPGPU este *CUDA (Compute Unified Device Architecture)*. Arhitectură dedicată procesării paralele, CUDA oferă programatorilor pasionați atât suportul hardware cât și uneltele software necesare implementării algoritmilor, dar și exemple și documentație, ce pot fi descărcate gratuit de pe site-ul producătorului [5].

Propunându-și a fi o continuare a abordării propuse în [6], această lucrare prezintă un studiu comparativ al performanțelor obținute în urma implementării unui algoritm de prelucrare a imaginilor în vederea execuției pe platforma CPU, respectiv pe platforma unui GPU din clasa GTX260, cu suport CUDA. Din volumul mare de rezultate obținute a fost selectat un eșantion considerat reprezentativ pentru ilustrarea situațiilor în care decizia de a paraleliza algoritmul și a-l rula pe platforma plăcii video conduce la performanțe deosebite, dar și a celor în care o astfel de variantă nu este adecvată.

## 2. Studiu de caz

Pentru a testa performanțele noii arhitecturi a fost proiectat un sistem software pentru calculul histogramei unei imagini pe 256 niveluri de gri. Acest tip de aplicație este o modalitate bună de test, având în vedere faptul că o imagine digitală este reprezentată ca o matrice de pixeli. Mai mult, numărul elementelor acestei matrice variază în funcție de dimensiunea imaginii, astfel încât pentru o imagine cu o rezoluție mare este necesar un volum de calcul substanțial. Histograma unei imagini, o unealtă foarte utilă în domeniul fotografiei digitale, are aspectul unui grafic ce indică poziția în care se găsesc toate nivelurile de luminozitate. Tonurile sunt dispuse pe orizontală, de la cele mai întunecate (în stânga) până la cele mai luminoase (în dreapta), iar cantitatea luminozității este reprezentată pe axa verticală. În total, histograma cuprinde 256 niveluri de luminozitate, de la 0 la 255 (de la negru pur la alb pur), care pot fi grupate în mai multe zone, în funcție de intensitate. Majoritatea aparatelor foto digitale, dar și mediile de editare precum Adobe Photoshop sau GIMP, pot afișa histograma unei imagini. Aplicația realizată cuprinde o bibliotecă statică pentru citirea secvențială (de pe disc) a imaginii în format BMP necomprimat, construirea unei matrice de pixeli și procesarea acesteia în vederea obținerii histogramei. Pentru o prezentare comparativă, algoritmul de construcție a histogramei a fost implementat atât în mod clasic, secvențial (și rulat pe CPU), cât și respectând considerentele de performanță specifice arhitecturii CUDA (și rulat pe placa video).

În ceea ce privește abordarea tradițională, în urma citirii pixelilor din imagine, au fost calculate valorile de intensitate luminoasă corespunzătoare fiecărui pixel, iar rezultatele au fost stocate în memoria centrală a sistemului. Pentru calculul histogramei, algoritmul a fost compilat și executat pe unitatea centrală de prelucrare (CPU). Pentru cea de-a doua abordare, respectiv paralelizarea algoritmului și rularea acestuia pe platforma procesorului grafic, s-a urmărit distribuția datelor de intrare și a calculelor pe mai multe fire de execuție. Astfel, fiecare subdiviziune, procesată de un *warp*, generează o sub-histogramă. În final, aceste componente sunt reunite pentru determinarea rezultatului final. La adaptarea metodei pe arhitectura procesorului grafic, trebuie avute în vedere o serie de constrângeri. Deși s-a urmărit o paralelizare a algoritmului cu o granularitate cât mai mare, trebuie precizat faptul că anumite operații (de exemplu, accesul la matricea de date inițială) rămân secvențiale. Astfel, execuția unui program tipic CUDA începe cu execuția codului specific CPU, iar atunci când acesta invocă un *kernel*, procesul de execuție este transferat procesorului grafic [3]. Conform recomandărilor prezentate în [6], pentru obținerea unor performanțe bune, s-a urmărit reducerea volumului de date stocate în memoria comună. Prin procesarea la nivel de *warp*, presiunea asupra memoriei comune a scăzut, utilizând din 16KB disponibili, numai 6KB ( $6 \text{ warps} \times 256 \text{ locații ale histogramei} \times 4 \text{ octeți/locație}$ ). Mai mult, grupând cele 192 fire de execuție ( $6 \text{ warps} \times 32 \text{ fire de execuție/warp}$ ) într-un bloc, au fost rulate două blocuri pe un multiprocesor, atingând astfel un grad ridicat de încărcare a GPU.

Performanțele aplicației sunt puternic influențate de configurația sistemului pe care rulează. Testul pentru acest studiu de caz a fost realizat pe un sistem dotat cu procesor AMD Athlon64 X2 la 2.2 GHz, 2 module de memorie RAM a câte 1 GB DDR400 care lucrează în dual-channel pe 128 biți (cu performanțe de scriere și de citire de 6300 MB/s la o latență de 54.2 ns), placă video Nvidia GTX260 ce oferă suport pentru arhitectura CUDA, cu 24 multiprocesoare (192 nuclee de procesare) și o frecvență GPU de 620 MHz. Prin intermediul unei interfețe pe 448 biți, procesorul grafic comunică cu memoria video GDDR3 (896 MB la o frecvență de 2,2 GHz, cu o lățime de bandă de 112 GB/s). Se poate observa diferența de performanță între cele două

cupluri de componente, CPU – memorie sistem și GPGPU - memorie video, de aproape 17,8 ori. Conexiunea între acestea este asigurată printr-o magistrală PCI Express 16x, iar lățimea de bandă dintre ele (conform unui test real al plăcii video) este 1300 MB/s, în comparație cu 4000 MB/s, valoare menționată în specificațiile producătorului.

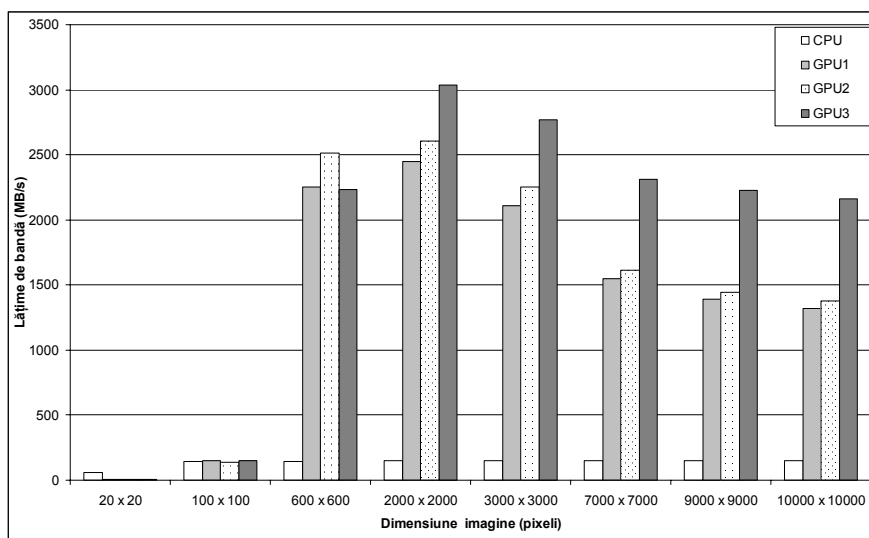
În tabelul 1.1 este prezentat un eșantion de date reprezentativ, obținut în urma testelor efectuate. Determinarea histogramei caracteristice unei imagini pe 256 tonuri de gri poate fi o operație costisitoare din punct de vedere al timpului de execuție atunci când algoritmul este conceput pentru a rula pe CPU, motiv pentru care am urmărit variația timpului de procesare ( $T_p$ ) și lățimea de bandă ( $B$ ) pentru dimensiuni diferite ale aceleași imagini, pentru fiecare unitate de prelucrare (procesor central, respectiv cele trei multiprocesoare de pe placa video). Se poate observa relativ ușor faptul că, pentru imagini mici, sub  $100 \times 100$  pixeli, unitatea centrală de procesare este mai rapidă decât procesorul grafic. În jurul acestei valori performanțele celor două procesoare sunt aproximativ identice, iar pentru dimensiuni mai mari se justifică paralelizarea algoritmului.

**Tabel 3.1. Rezultatele obținute pentru dimensiuni diferite ale imaginii**

Dimensiune imagine (pixeli)	CPU		GPU1		GPU2		GPU3	
	$T_p$ (ms)	$B$ (MB/s)	$T_p$ (ms)	$B$ (MB/s)	$T_p$ (ms)	$B$ (MB/s)	$T_p$ (ms)	$B$ (MB/s)
20 x 20	0,006	56,895	0,074	5,090	0,072	5,249	0,064	5,892
100 x 100	0,065	144,649	0,064	146,925	0,068	139,691	0,063	149,971
600 x 600	2,405	142,717	0,152	2254,805	0,136	2515,901	0,153	2232,025
2000 x 2000	25,4	150,184	1,559	2445,454	1,462	2608,382	1,256	3036,442
3000 x 3000	56,887	150,609	4,072	2107,800	3,804	2255,750	3,099	2769,310
7000 x 7000	310,575	150,462	30,223	1546,158	28,989	1611,979	20,231	2309,738
9000 x 9000	513,216	150,516	55,639	1388,354	53,469	1444,717	34,684	2227,158
10000 x 10000	634,374	150,333	72,305	1318,957	63,318	1375,795	44,113	2161,855

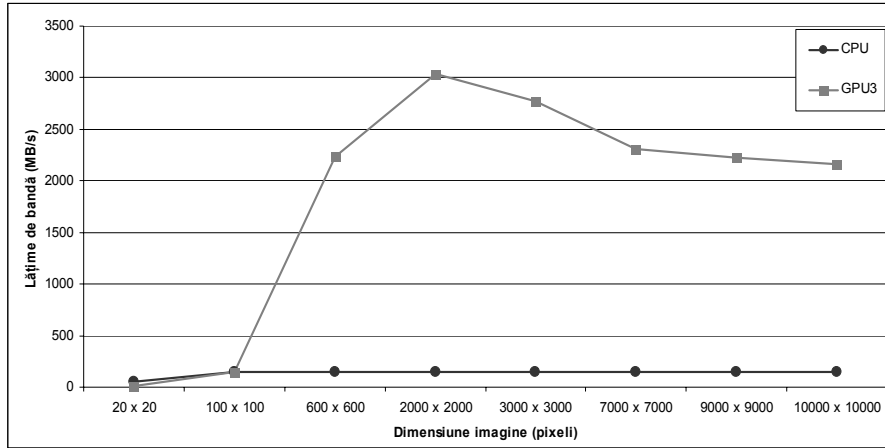
Reprezentarea grafică a rezultatelor obținute ca urmare a testelor efectuate evidențiază și mai clar performanțele specifice celor două abordări. Întrucât plaja de valori a rezultatelor este destul de cuprinzătoare, sunt prezentate separat cazul imaginilor de dimensiuni mici, respectiv al celor de dimensiuni mari.

Analizând informațiile reprezentate în figura 3.1 se poate observa că lățimea de bandă pentru imagini de dimensiuni mari este influențată în mod direct de dimensiunea imaginii. De asemenea, se constată că reducerea acesteia este mai mare în cazul variantei fără operații atomice la nivel de memorie.



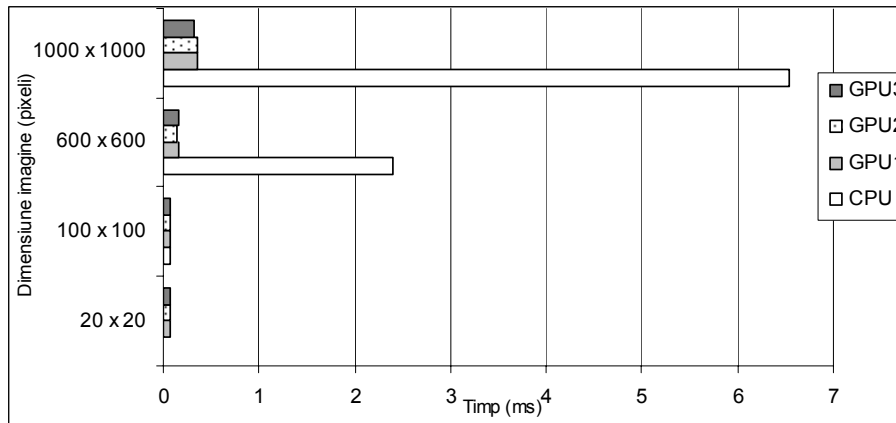
**Figura 3.1. Variația lățimii de bandă în funcție de dimensiunea imaginii**

Ilustrată grafic în figura 3.2, variația lățimii de bandă apărută la nivelul procesorului central, respectiv la nivelul procesorului grafic subliniază faptul că, pentru imagini relativ mici, unitatea centrală de procesare este totuși mai rapidă decât procesorul grafic. Cu toate acestea, începând cu imagini de dimensiuni mai mari de  $100 \times 100$  pixeli, performanța obținută în urma procesării pe platforma grafică este semnificativ mai bună, atingând o valoare maximă pentru imagini de  $2000 \times 2000$  pixeli.

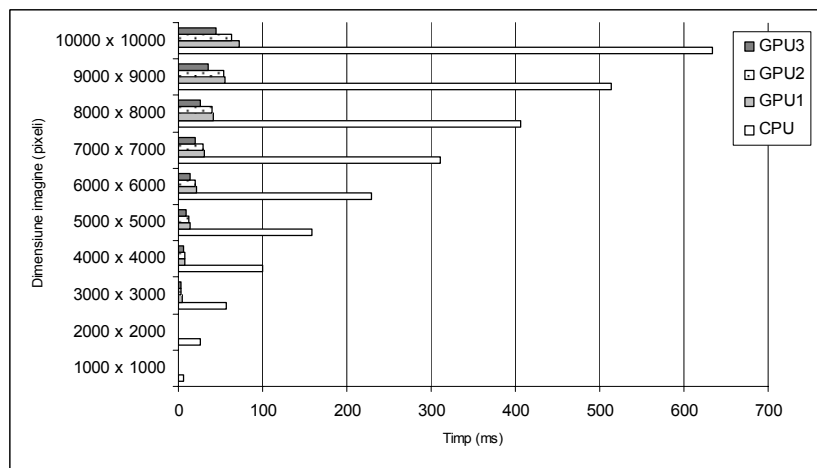


**Figura 3.2. Comparație între variațiile lățimii de bandă pentru CPU și GPU**

După cum s-a precizat anterior, dimensiunea imaginii reprezintă un factor determinant în variația duratei de execuție a algoritmului de prelucrare. Pe baza rezultatelor obținute în urma testării, în figurile 3.3 și 3.4 au fost reprezentate variațiile duratei de execuție pentru prelucrarea imaginilor de dimensiuni mici și respectiv mari.

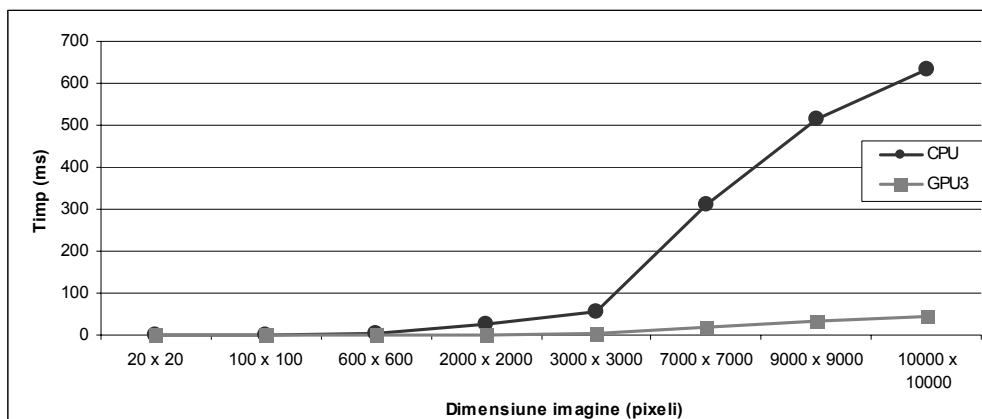


**Figura 3.3. Variația duratei de execuție pentru imagini cu dimensiuni mici**



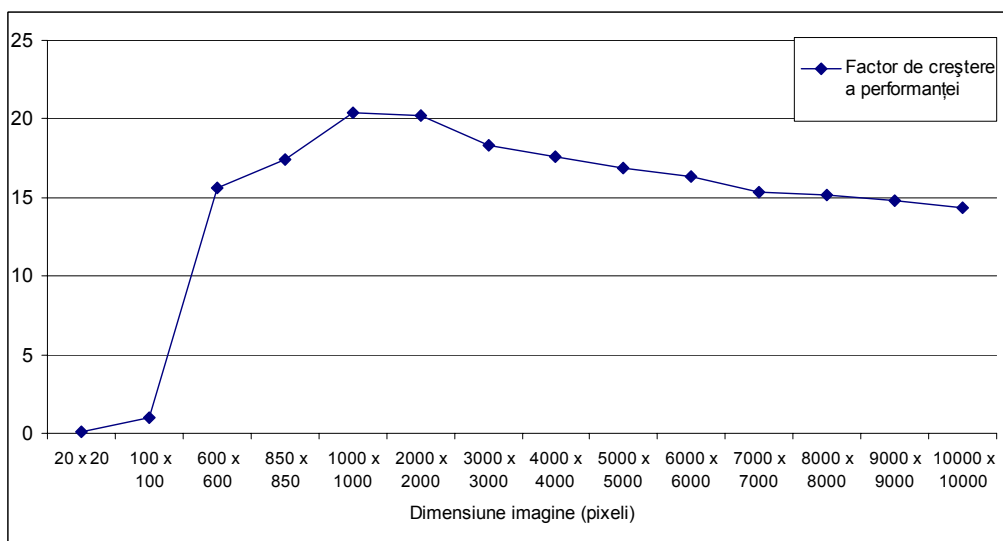
**Figura 3.4. Variația duratei de execuție pentru imagini cu dimensiuni mari**

Utilizând același eșantion de date prezentat anterior în tabelul 1.1, figura 3.5 ilustrează grafic o analiză comparativă a rezultatelor obținute, cu referire la variația duratelor de execuție în cazul rulării algoritmului pe platforma CPU, respectiv pe platforma GPU.



**Figura 3.5. Variația duratei de execuție în raport cu dimensiunea imaginii**

Pentru valorile obținute în cadrul studiului de caz, s-a considerat ca relevantă calcularea unui factor de creștere a performanțelor definit ca raportul dintre durata de execuție a algoritmului pe CPU și durata de execuție pe GPU, în cazul utilizării memoriei comune. Reprezentarea grafică a acestui factor în funcție de dimensiunea imaginii se regăsește în figura 3.6.



**Figura 3.6. Reprezentarea factorului de creștere a performanței**

Rezultatele studiului de caz conduc la concluzia că performanțele obținute prin utilizarea arhitecturii CUDA sunt semnificativ mai bune în comparație cu cele obținute prin rularea aplicației pe platforma unității centrale de procesare. În cadrul studiului s-a constatat o eficiență mai mare a arhitecturii pentru imagini cu dimensiuni cuprinse între 2 și 9 megapixeli, interval pentru care a fost obținută și cea mai mare lățime de bandă, aproximativ 3000 MB/s. Totuși, se poate observa că pentru imagini de dimensiuni mici efortul de paralelizare nu este justificat, performanțele obținute în urma execuției variantei tradiționale (pe platforma CPU) fiind aproximativ aceleași.

### 3. Concluzii

Scopul principal al acestei lucrări a fost cercetarea potențialului de performanță promis de utilizarea arhitecturii CUDA pentru calcul paralel. Rezultatele studiului de caz au confirmat faptul că utilizarea acestei tehnologii pentru rezolvarea problemei de calculare a histogramei

unei imagini pe 256 tonuri de gri poate conduce la obținerea unor performanțe deosebite. Chiar și în urma optimizării codului sursă, rularea algoritmului pe unitatea centrală de prelucrare s-a dovedit a fi o operație costisitoare din punctul de vedere al duratei de execuție. Utilizarea arhitecturii CUDA a condus la obținerea unui factor de optimizare cu ordin de multiplicare mai mare de 20, pentru imagini cuprinse între 1 și 4 mega-pixeli. Cu toate acestea, trebuie menționat faptul că, în cazul procesării unui volum mic de date, utilizarea algoritmilor de calcul paralel nu este justificată. Dimensiunea de  $100 \times 100$  pixeli, pragul peste care performanțele obținute pentru procesarea pe GPU au fost semnificativ mai bune decât cele pentru procesarea pe CPU, este mică față de dimensiunile imaginilor folosite în mod curent.

Așadar, în acest caz, tehnologia CUDA și-a dovedit eficiența, chiar dacă implementarea algoritmului nu a fost concentrată pe optimizarea codului, respectiv pe managementul direct al resurselor puse la dispoziție de NVIDIA. Acest rezultat poate fi susținut și de faptul că, în ultimul timp, multe companii cu nume sonor și-au îndreptat atenția către arhitectura paralelă, cu nuclee multiple, a unităților de procesare grafică Nvidia[5]. Domenii precum procesarea audio-video, modelarea și simularea în timp real a proceselor cu grad ridicat de complexitate, proiectarea electronică automată sau prelucrarea imaginilor din diverse domenii, toate pot avea în comun suportul oferit de tehnologia CUDA pentru dezvoltarea și utilizarea unor algoritmi de o complexitate mult mai mare decât ar fi fost posibil altfel și o durată de execuție de zeci, chiar sute de ori mai mic. Cu toate acestea, CUDA nu este singura opțiune pentru implementarea algoritmilor de calcul și exploatarea platformei procesorului grafic. În multe privințe, platforma celor de la RapidMind de exemplu, *Multicore Development Platform* (achiziționată de către Intel în 2009), este superioară platformei CUDA [2]. Totuși, ca și în trecut, concurența de pe piața microprocesoarelor nu poate decât să susțină evoluția tehnologică din domeniul IT.

## BIBLIOGRAFIE

1. **CASTANO-DIEZ, D.; MOSER D.; SCHOENEGGER A.; PRUGGNALLER S.; FRANGAKIS A. S.** Performance Evaluation of Image Processing Algorithms on the GPU. *Journal of Structural Biology*, Vol. 164, 2008, pp. 153 – 160.
2. **HALFHILL, T.** Parallel Processing with CUDA – Nvidia’s High – Performance Computing Platform Uses Massive Multithreading, 2008, documentație Internet, [http://www.nvidia.com/docs/IO/55972/220401\\_Reprint.pdf](http://www.nvidia.com/docs/IO/55972/220401_Reprint.pdf)
3. **KIRK, D.; HWU W.** Programming Massively Parallel Processors, Curs Universitatea din Illinois, 2007.
4. \*\*\* General-Purpose Computation on Graphics Hardware, documentație Internet, <http://ggpu.org/>
5. \*\*\* Nvidia Documentation, 2009, documentație Internet, <http://www.nvidia.com>
6. **MOCANU, Ș.; DOBRESCU R.; SARU D.; DIN R.; GRUMĂZESCU A.** Arhitecturi Complexe folosite în prelucrarea paralelă a imaginilor. *Revista Română de Automatică și Informatică*, vol. 20, 2010, pp. 97 – 105.