

Alexandru PÎRJAN

**SOLUȚII INOVATOARE DESTINATE
OPTIMIZĂRII PROCESĂRII
DATELOR**



Copyright © 2013, **Editura Pro Universitaria**

Toate drepturile asupra prezentei ediții aparțin
Editurii Pro Universitaria

Nicio parte din acest volum nu poate fi copiată fără acordul scris al
Editurii Pro Universitaria

Descrierea CIP a Bibliotecii Naționale a României

PÎRJAN, ALEXANDRU

**Soluții inovatoare destinate optimizării procesării
datelor / Alexandru Pîrjan. - București : Pro Universitaria,
2013**

Bibliogr.

ISBN 978-606-647-611-9

004(075.8)

PREFAȚĂ

Conținutul științific al acestei cărți reprezintă o parte din rezultatele obținute în urma cercetărilor desfășurate pe parcursul programului de studii doctorale din cadrul Academiei de Studii Economice din București, sub îndrumarea științifică a domnului Prof. univ. dr. Ion Lungu [PiTe12]. Teza de doctorat „Soluții de optimizare a procesării datelor” a fost susținută pe data de 12 septembrie 2012, la Academia de Studii Economice din București, iar ulterior, prin Ordin al ministrului educației, cercetării, tineretului și sportului, autorului i-a fost atribuit titlul de doctor în domeniul Informatică Economică.

Pentru cercetarea stadiului actual al cunoașterii în domeniul procesării datelor au fost utilizate o serie de publicații științifice consacrate: reviste cotate ISI, reviste de circulație internațională, standarde, metodologii consacrate, cărți publicate de edituri internaționale recunoscute (Springer, Elsevier), cât și volumele conferințelor naționale și internaționale de profil.

Scopul cercetării a constat în identificarea de soluții pentru îmbunătățirea performanței funcțiilor algoritmice de bază ce intervin frecvent în pașii de calcul ai unui număr mare de algoritmi sau aplicații de procesare a datelor, a căror optimizare implică resurse computaționale și economice semnificative. Odată identificate funcțiile, a fost studiată posibilitatea implementării acestora în arhitectura de procesare paralelă CUDA (Compute Unified Device Architecture) pentru a beneficia de puterea computațională imensă oferită de aceasta. Au fost dezvoltate o serie de soluții de optimizare a performanței funcțiilor algoritmice de bază în CUDA, construind astfel o bibliotecă de funcții algoritmice de bază (BFAB) [PiTe12], ce facilitează dezvoltarea optimă a algoritmilor și a aplicațiilor, implicit optimizarea procesării datelor în ansamblu. În consecință, procesarea datelor a fost optimizată prin dezvoltarea soluției BFAB ce conține funcțiile algoritmice, structurate ca blocuri funcționale de bază, ce intră în componența unui număr mare de algoritmi și aplicații de procesare a datelor. În această carte sunt prezentate doar o parte dintre funcțiile algoritmice de bază din BFAB și dintre aplicațiile lor, dezvoltate și prezentate în cadrul tezei de doctorat a autorului [PiTe12].

Motivația cercetării este dată de actualitatea și importanța temei abordate, ce rezultă din modul cum a evoluat tehnologia modernă privind procesarea datelor. Una dintre problemele curente este absența unor metode performante și accesibile din punct de vedere economic care să permită implementarea și dezvoltarea de soluții de optimizare a procesării datelor bazate pe cele mai noi arhitecturi hardware și software.

În ultimii ani, creșterea fără precedent a volumului de date ce trebuie procesat, cât și dezvoltarea aplicațiilor de procesare de date în diverse domenii (comercial, industrial, administrativ, etc) au determinat intensificarea interesului cercetătorilor cu privire la dezvoltarea soluțiilor de optimizare a procesării datelor. Un factor determinant ce a dus la creșterea impresionantă a volumului de date într-o varietate de domenii economice și științifice, îl reprezintă evoluția continuă și rapidă a tehnicii de calcul, a dispozitivelor și mediilor de stocare cât mai ales a arhitecturilor de procesare paralelă. Pentru procesarea acestor volume imense de date sunt necesare resurse costisitoare și greu accesibile. În majoritatea situațiilor este esențial ca durata de timp necesară procesării datelor să fie cât mai mică.

Procesarea eficientă a datelor este extrem de utilă în numeroase domenii de activitate: în domeniul afacerilor, în predicția datelor financiare, în telecomunicații, în neuroștiințe, în analiza datelor medicale, în explorarea resurselor naturale, în crearea de aplicații inovatoare în domenii precum recunoașterea imaginilor, redarea grafică de înaltă definiție în timp real și codificare. Astfel, pentru a planifica o strategie economică eficientă, se impune utilizarea unor tehnici adecvate pentru a procesa și analiza datele disponibile. Cele mai multe companii utilizează volume mari de date în activitățile lor economice și în consecință necesită soluții eficiente de procesare a datelor. Pentru a supraviețui, pentru a se dezvolta și a obține profit în domeniul competitiv al afacerilor, companiile trebuie să gestioneze și să proceseze cât mai eficient posibil datele disponibile, obținând astfel strategii economice adaptate nevoilor și cerințelor pieței.

Eficiențizarea componentei software și utilizarea optimă a resurselor hardware disponibile sunt cei mai importanți factori în dezvoltarea soluțiilor eficiente de procesare a datelor, în identificarea modelelor și în extragerea interdependențelor dintre diferitele categorii de date. Aceste metode ajută la prezicerea viitoarelor modele ale tendințelor pieței și oferă astfel companiilor ce le folosesc un avantaj competitiv. De exemplu: datele burselor de valori

sunt procesate în vederea elaborării de predicții; motorul de căutare Google utilizează tehnici de procesare a datelor pentru a selecta, din milioanele de pagini disponibile pe Internet, paginile Web cele mai relevante pentru căutarea cerută de client; în vederea stabilirii locației, datele geospațiale necesită procesare și interpretare.

În privința proiectelor de cercetare științifică, este important ca analiza datelor experimentale să valideze ipotezele științifice pentru care sunt proiectate experimentele. Un exemplu de demers științific actual ce folosește procesarea datelor este acceleratorul de particule CERN de lângă Geneva, "Large Hadron Collider" (LHC), care generează zilnic o cantitate imensă de date, a căror analiză și procesare este realizată automat, prin utilizarea unor aplicații software și a unor sisteme hardware performante.

Progresele înregistrate de tehnologia informației în ultimii ani au determinat companiile să acumuleze volume imense de date din diverse domenii științifice sau economice. Costurile ce decurg din stocarea acestor volume de date au impus dezvoltarea corespunzătoare a tehnicii de calcul, a unor dispozitive și medii de stocare, a unor instrumente noi care să permită și să faciliteze atât accesul la date cât și procesarea acestora.

Deși puterea de calcul evoluează odată cu dezvoltarea noilor tehnici și instrumente, volumul de date crește din ce în ce mai mult și se constată că timpul alocat procesării datelor are, de asemenea, tendința să crească. În vederea reducerii costurilor și a duratei de timp necesare procesării datelor, se impune dezvoltarea unor metode avansate și performante de procesare a datelor pentru a gestiona în mod corespunzător volumul imens de date, precum și a unor soluții de optimizare a procesării datelor, folosind atât funcții algoritmice de bază eficiente cât și arhitecturi hardware performante.

Pe baza problemelor identificate, în cadrul cercetării au fost atinse următoarele obiective principale:

- identificarea funcțiilor algoritmice cu cele mai bune performanțe, pentru a fi utilizate în dezvoltarea aplicațiilor privind procesarea datelor;
- cercetarea arhitecturilor hardware și a mediilor de programare care să ofere soluțiile optime de implementare a funcțiilor algoritmice de bază;
- dezvoltarea de soluții pentru îmbunătățirea performanței funcțiilor algoritmice de bază ce intervin frecvent în pașii de calcul ai unui număr mare de algoritmi sau aplicații de procesare a datelor. Aceste soluții au fost proiectate astfel încât să permită obținerea

unor rezultate net superioare din punct de vedere al performanței și costurilor față de cele ale metodelor utilizate până în prezent, implementate în cadrul arhitecturilor clasice care folosesc unitățile de procesare centrală.

Pe baza concluziilor rezultate din experimentele efectuate, au fost identificate soluții de optimizare a funcțiilor algoritmice de bază prin adaptarea acestora în vederea executării în paralel a anumitor pași de execuție, paralelizarea execuției fiind favorizată și recomandabilă întrucât rularea funcțiilor algoritmice are loc pe unități de procesare grafică, ce implementează arhitectura CUDA (Compute Unified Device Architecture).

În urma testărilor efectuate, a rulărilor experimentale pe procesoare grafice, s-au obținut rezultate deosebit de promițătoare, ceea ce conferă soluțiilor dezvoltate un înalt nivel de performanță și utilitate. Soluțiile au fost proiectate astfel încât să fie auto-ajustabile, să aloce resursele în funcție de arhitectura de procesare folosită (Tesla, Fermi sau Kepler), ceea ce le conferă un grad ridicat de originalitate. Pe baza volumului mare de date de intrare generate, a unei game variate de arhitecturi, a numărului mare de iterații ale testelor efectuate, s-a obținut o analiză detaliată a caracteristicilor funcțiilor algoritmice de bază folosite în optimizarea procesării datelor.

Testele au fost rulate pe următoarea configurație: Intel i7-2600K la 4,6 GHz cu 8 GB (2x4GB) de memorie cu frecvența 1333 MHz, DDR3, rulând în modul dual channel. Unitățile de procesare grafică utilizate au fost GeForce GTX 280, GTX 480 și GTX 680. Pentru a programa și accesa unitatea de procesare grafică s-a utilizat CUDA toolkit 4.1, versiunile de driver Nvidia 270.81 pentru GTX 280 și pentru GTX 480 și versiunea Nvidia 301.10 pentru GTX 680. Sistemul de operare utilizat în cadrul testelor experimentale a fost Windows 7 pe 64 biți. În plus, toate procesele legate de interfața grafică a utilizatorului au fost dezactivate pentru a reduce traficul extern al unității de procesare grafică. În măsurătorile efectuate nu au fost incluși timpii necesari transferurilor de date între unitatea de procesare centrală (CPU) și unitatea de procesare grafică (GPU), întrucât funcțiile algoritmice de bază au fost proiectate pentru a fi folosite în componența unui număr mare de aplicații ce rulează pe procesoare grafice, astfel încât timpii de transfer variază în funcție de complexitatea și specificul aplicației. Pentru a obține o evaluare fiabilă a timpului de execuție necesar rulării funcțiilor algoritmice de bază pe unitățile de procesare grafică, au fost folosite evenimentele de înregistrare a ștampilelor de timp din cadrul interfeței de programare a aplicației (API) CUDA.

Experimentele au fost efectuate pe cele mai noi și reprezentative platforme hardware cu diferite arhitecturi și configurații. Rezultatele experimentale obținute sunt prezentate în această lucrare atât sub formă tabelară, cât și sub formă grafică, împreună cu o serie de observații aferente acestor seturi de teste. Prin intermediul rulărilor experimentale au fost măsurate și analizate atât timpul mediu de rulare cât și lățimile de bandă, folosind 10.000 de iterații pentru fiecare test al funcțiilor algoritmice.

În urma analizării rezultatelor cercetării prezentate în acest material, se obține o imagine de ansamblu a performanțelor oferite de funcțiile algoritmice de bază implementate pe unitățile de procesare centrală dar și a soluțiilor de optimizare a aplicațiilor dezvoltate în CUDA. Analizând această imagine, se ajunge la concluzia că o îmbunătățire substanțială a timpilor de execuție aferenți procesării datelor se obține prin rularea funcțiilor algoritmice de bază pe procesoare suplimentare, care conlucrează cu procesorul dedicat în operațiile de procesare a datelor. Pentru atingerea acestui scop există diverse soluții, ca de exemplu:

- utilizarea mai multor procesoare centrale CPU, ceea ce implică costuri suplimentare considerabile cu privire la platforma hardware, necesitând modificări substanțiale, cum ar fi placa de bază multi-soclu, etc.
- folosirea unei unități de procesare grafică de scop general (GPGPU) ceea ce nu necesită modificări substanțiale ale arhitecturii hardware folosite (placa de bază, memorii, etc).

Din motive economice a fost preferată a doua variantă, iar rezultatele experimentale obținute au confirmat faptul că aceasta reprezintă o soluție performantă. Pentru a optimiza funcțiile algoritmice de bază folosite în procesarea datelor, s-a pornit de la principiul utilizării unităților de procesare grafică cu mai multe nuclee de procesare care să suporte ca mediu de dezvoltare arhitectura paralelă CUDA (Compute Unified Device Architecture). Această arhitectură permite utilizatorilor să folosească limbajul CUDA C, iar din punct de vedere operațional, procesorul grafic devine un coprocesor al unității de procesare centrală. În vederea procesării în paralel a codului sursă, au fost identificate toate operațiile ce pot fi paralelizate în cadrul arhitecturii CUDA. De asemenea, a fost acordată o atenție deosebită modului de acces la diversele tipuri de memorie: privată, partajată sau globală, luând în considerare caracteristicile și timpii de acces specifici acestora. Tehnicile de proiectare utilizate pe parcursul cercetării, în

contextul rulării pe unități de procesare grafică, sunt aplicabile pe o gamă variată de procesoare ce încorporează mai multe nuclee de execuție.

Cartea este structurată în 6 capitole, cuprinzând de asemenea un capitol de concluzii, bibliografie și anexe. În dezvoltarea cercetării s-a urmărit o abordare graduală ce corespunde următoarei structuri pe capitole:

După ce în **Prefață** este prezentată motivația, scopul și conținutul cercetării, în **Capitolul 1** este descrisă arhitectura CUDA (Compute Unified Device Architecture) de optimizare a procesării datelor, pornind de la un scurt istoric al evoluției de la unități de procesare centrală, la unități de procesare grafică. Sunt prezentate apoi câteva elemente referitoare la CUDA C, primul limbaj proiectat special de o companie producătoare de procesoare grafice pentru a facilita calculele cu scop general pe unitățile de procesare grafică. Sunt descrise o serie de domenii în care au fost implementate cu succes limbajul CUDA C și arhitectura CUDA, precum și cele trei niveluri de abstractizare oferite de mediul de programare CUDA și componentele de bază ale arhitecturii CUDA.

Sunt prezentate apoi o serie de concepte de bază ale modelului de programare CUDA: ierarhia firelor de execuție, ierarhia memoriei, echipamentul hardware gazdă și cel specializat. De un deosebit interes în cadrul acestui demers științific a fost comparația celor mai importante arhitecturi de procesoare grafice ce implementează CUDA. De asemenea sunt analizate principalele avantaje și limitări ale utilizării arhitecturii CUDA în dezvoltarea aplicațiilor software, oportunitatea dezvoltării unei soluții de optimizare a procesării datelor în CUDA, iar apoi sunt identificate o serie de domenii cheie în care au fost implementate cu succes limbajul CUDA C și arhitectura CUDA.

Aspectul esențial din cadrul acestui capitol se referă la identificarea celei mai importante facilități oferite de arhitectura CUDA: pentru a beneficia de puterea de procesare a unităților GPU, dezvoltatorii au posibilitatea să scrie aplicațiile în diverse limbaje de programare, iar pentru dezvoltarea cu succes a aplicațiilor performante, programatorii trebuie să aiba la dispoziție biblioteci specializate de funcții puternice dezvoltate în CUDA, ce pot fi apelate cu ușurință în cadrul aplicațiilor și au ca efect creșterea substanțială a performanței acestora. Aceasta a fost și direcția de cercetare abordată în cadrul acestui demers științific, anume dezvoltarea unei soluții de optimizare a performanței funcțiilor algoritmice de bază, construind astfel o bibliotecă de funcții algoritmice de bază (BFAB) [PiTe12] în CUDA, apelabile în

programarea aplicațiilor de procesare de date. În finalul capitolului sunt prezentate oportunitatea, utilitatea și importanța dezvoltării soluției BFAB în procesarea datelor, evidențiind interfața de programare (API) a acesteia. În această carte sunt prezentate următoarele funcții algoritmice de bază, din cadrul soluției BFAB [PiTe12]: însumarea paralelă prefixată, reducerea paralelă și compactarea fluxurilor de date.

În **Capitolul 2** sunt prezentate soluții de îmbunătățire a performanței funcției algoritmice de însumare paralelă prefixată în CUDA. Este prezentată mai întâi proiectarea acestei funcții în cadrul warp-urilor, apoi în cadrul blocurilor de fire de execuție, iar în final la nivel global. A fost dezvoltată o variantă optimizată a funcției algoritmice de însumare prefixată în CUDA, au fost identificate și aplicate o serie de soluții de îmbunătățire a performanței acesteia. Analizând și interpretând rezultatele experimentale obținute, se constată că acestea confirmă eficiența funcției algoritmice de însumare prefixată din BFAB, ce oferă rezultate optime în diverse situații și astfel, acesta este implementabilă într-o gamă largă de algoritmi și aplicații de procesare de date.

În **Capitolul 3** sunt prezentate soluții de îmbunătățire a performanței funcției algoritmice de reducere paralelă în CUDA. Este prezentată funcția algoritmică de reducere paralelă în CUDA și variante de optimizare ale acesteia printr-o serie de soluții de îmbunătățire a performanței. Eficiența funcției algoritmice de reducere paralelă din cadrul soluției BFAB a fost confirmată prin analiza și interpretarea rezultatelor experimentale obținute. Se constată că funcția algoritmică de reducere paralelă oferă rezultate optime în diverse situații și, întrucât aceasta nu depinde de operatorul binar utilizat, poate fi implementată într-o gamă largă de algoritmi și aplicații.

În **Capitolul 4** sunt prezentate soluții de îmbunătățire a performanței funcției algoritmice de compactare paralelă a fluxurilor de date în CUDA. Pentru început este descrisă proiectarea acesteia, apoi este dezvoltată o variantă optimizată a funcției algoritmice de compactare a fluxurilor de date în CUDA, utilizând o serie de soluții de îmbunătățire a performanței acesteia. Rezultatele experimentale obținute confirmă eficiența funcției de compactare din BFAB, ce oferă rezultate optime în diverse situații și, astfel, funcția este implementabilă într-o gamă largă de algoritmi și aplicații de procesare de date.

În cadrul **Capitolului 5** sunt analizate o serie de considerații economice privind oportunitatea soluțiilor de optimizare a procesării datelor

în CUDA. Pentru început, sunt evidențiate o serie de avantaje economice ale utilizării GPU legate de îmbunătățirea considerabilă a eficienței energetice în industria de calculatoare. Sunt prezentate apoi câteva dintre domeniile în care utilizarea sistemelor bazate pe GPU a adus o reducere semnificativă a timpului de lucru și a cheltuielilor aferente consumului de energie. Această analiză reliefează importanța din ce în ce mai mare a unităților de procesare grafică în întreaga industrie de calcul, ajungând astfel la concluzia că GPU reprezintă soluții de calcul eficiente din punct de vedere economic. În continuare este prezentat un studiu economic privind alegerea tehnologiei CUDA în implementarea soluției BFAB, punând în evidență, pentru început, o serie de aspecte ce trebuie luate în considerare atunci când se evaluează aspecte economice referitoare la utilizarea procesoarelor grafice ce implementează tehnologia CUDA. Pe baza acestor aspecte se remarcă oportunitatea soluțiilor de optimizare a procesării datelor în CUDA. În final este prezentată o analiză referitoare la avantajele economice ale soluției de optimizare a procesării datelor în CUDA, luând în considerare costul consumului de energie electrică și cel al echipamentului hardware. Rezultatele experimentale confirmă avantajele economice ale soluției propuse.

În **Capitolul 6** este propusă o soluție software de lansare, rulare și testare experimentală a funcțiilor algoritmice de bază din BFAB, dezvoltată în Windows Presentation Foundation (WPF), folosind limbajul declarativ Extensible Application Markup Language (XAML) și limbajul de programare C#. Folosind documentația oficială Microsoft, precum și literatura de specialitate, pentru început este analizat și prezentat modelul de programare WPF, principalele sale caracteristici, avantajele separării între interfața grafică și funcționalitatea aplicației, tipurile de aplicații ce pot fi dezvoltate în WPF, tipurile de navigare suportate de modelul de programare WPF, arhitectura WPF, componentele majore ale acesteia, ierarhia claselor frecvent utilizate în WPF, posibilitatea de legare la sursele de date, o serie de caracteristici și facilități ce au fost apoi implementate în dezvoltarea soluției software de lansare, rulare și testare experimentală a funcțiilor algoritmice de bază din BFAB. Sunt prezentate apoi elemente legate de proiectarea acestei soluții, de validarea datelor introduse, precum și un exemplu de rulare a soluției dezvoltate în WPF. Soluția dezvoltată în cadrul acestui capitol are un grad ridicat de originalitate și este utilă în lansarea, testarea și rularea funcțiilor algoritmice de bază prezentate și analizate în capitolele anterioare.