

**Advanced research methods for multi-objective optimization  
of CPUs architectures**

**Metode avansate de cercetare pentru optimizarea multi-  
obiectiv a arhitecturilor de procesare**

Teză de abilitare

**Rezumat**

Adrian FLOREA

2021

Această teză de abilitare reprezintă un sumar al celor mai importante realizări personale în domeniul arhitecturilor de procesoare, cu accent pe perioada ulterioară susținerii tezei de doctorat. Evident că, cu mici excepții, la munca din spatele obținerii acestor realizări am avut cel puțin un colaborator, fie pe mentorul științific – dl. profesor Vișan Lucian, fie pe alți colegi din universitate sau din străinătate (co-autori la lucrările științifice publicate). Astfel, este oarecum nedrept că există un singur autor al tezei de abilitare deoarece datorez recunoștință colaboratorilor mei pentru tot ce am realizat. Pe de altă parte, cel puțin în domeniul microarhitecturilor, sunt foarte puține cazuri în care activitatea de cercetare este depusă de o singură persoană, de cele mai multe ori sunt colective de 3-4 cercetători care lucrează împreună pe anumite tematici. În afara acestor contribuții, intenția este de a ilustra viziunea mea asupra domeniului arhitecturilor de procesoare în care sunt plasate contribuțiile mele anterioare. Având în vedere că cercetările mele în acest vast și de actualitate domeniu datează de mai bine de douăzecișidoi de ani, această sinteză a activității din ultimii șaisprezece ani se întinde pe un spectru destul de larg ținând cont și de dinamicitatea domeniului, rezultatele fiind inevitabil legate de efortul depus în întreaga activitate de cercetare. Pe lângă caracterul științific, lucrarea are un puternic aspect tehnic datorită multor simulări dezvoltate. Autorul sintetizează metodologiile de simulare utilizate în cercetare, instrumentele software de compilare, depanare și simulare, benchmark-urile standardizate SPEC (*Standard Performance Evaluation Corporation*) pentru arhitecturile monocore și SPLASH (*Stanford Parallel Applications for Shared memory*) pentru arhitecturile multicore. Deși valoarea adăugată este mai degrabă de natură tehnică și nu științifică, utilitatea sa în cadrul tezei de abilitare este indiscutabilă.

Primul capitol cuprinde o motivație a implicării mele în acest domeniu, o introducere în problematica microarhitecturilor provocările și limitările existente, precum și calea pe care am urmat-o în a sintetiza prezenta lucrare. Totodată, am încercat să punctez elementele importante ale activității mele de cercetare și să rezum modestele mele contribuții la dezvoltarea domeniului. Impactul cercetărilor a fost reliefat în indici scientometrici și în prestigiul publicațiilor ce citează lucrările publicate. Un lucru demn de subliniat este faptul că peste 42% din articolele publicate au ca și co-autor un student la una din formele de învățământ (licență, master sau doctorat). Colaborarea a fost extrem de utilă pentru ambele părți profesor / student în vederea inițierii studenților în cercetare și creării unei structuri în permanentă dezvoltare la nivelul centrului de cercetare Advanced Computer Architecture and Processing Systems (ACAPS) (<http://acaps.ulbsibiu.ro/index.php>) prin care să aibă loc transferul de cunoștințe de la profesor la doctoranzi, de la doctoranzi la masteranzi și de la masteranzi la studenți, fiecare dintre aceștia parcurgând treptele dezvoltării profesionale și învățând de la înaintașii (mentorii) lor. Personal am muncit alături de ei și am avut împreună satisfacția muncii împlinite.

Capitolul al doilea prezintă limitarea *Memory Wall* cauzată de creșterea decalajului dintre viteza procesorului și lentoarea memoriei principale (DRAM) și necesitatea folosirii unui sistem ierarhic bazat pe unul sau mai multe niveluri de memorii cache, separate sau unificate. Contribuțiile personale au fost în ton cu dezvoltările tehnologice ale vremurilor și provocările curente ale domeniului, de la concepte precum (Selective) Victim Cache (propușe inițial de Norman Jouppi [1], continuate de Stiliadis [2], implementate în procesoarele comerciale AMD [3] și folosit ca și concept fundamental în alte investigații de tipul *cache dead blocks* [4], *victim buffer in embedded*

*system design* [5], etc) la influența negativă în procesoare cu mii de instrucțiuni *in-flight* produsă de aplicațiile de tip *pointer chasing* în care are loc serializarea execuției operațiunilor cu referire la memorie caracterizate de latență ridicată. Abordările recente specifice memoriilor cache vizează cercetarea acestora nu doar din perspectivă tehnică ci și impactul la nivel societal actual, în principal, cerința de a reduce consumul de energie și de a reduce amprenta de carbon a algoritmilor de înlocuire a blocurilor din cache precum și asigurarea securității datelor memorate.

O provocare continuă în știința și ingineria calculatoarelor o reprezintă înțelegerea profundă a complexelor și subtilelor interacțiuni hardware-software. Optimizările interfeței hardware-software nu sunt posibile decât pe baza unui proces de înțelegere, care necesită o viziune integratoare asupra multiplelor și complicatele procese de prelucrare a informației, codificată la niveluri semantice diferite. În acest sens, în capitolul al treilea sunt prezentate cercetările autorului relative la analiza legăturilor adânci dintre caracteristicile unor programe procedurale și obiectuale HLL (High Level Languages), pe de o parte, și execuția lor pe arhitecturile cu paralelism la nivelul instrucțiunilor, respectiv structurile hardware de predicție a salturilor indirecte prin registru, pe de altă parte. O parte din contribuțiile ilustrate în acest capitol au fost incluse în teza de doctorat, etapă fundamentală în cariera mea profesională, care a culminat cu formarea mea ca cercetător. Aportul adus relativ recent referitor la influența compilatorului asupra arhitecturilor de procesare s-a concretizat prin implementarea unor metode de optimizare bazate pe algoritmi genetici și de inteligență colectivă (*swarm intelligence*) pe problema de colorare a grafurilor aplicată pentru alocarea regiștrilor de către compilator în sistemele încorporate.

În capitolul 4 este prezentată limitarea fundamentală a paradigmei *ILP Wall* intitulată *fetch bottleneck / control flow bottleneck* (limitarea producătorului) și metodele predictive de creștere a performanței aplicate fluxului de control al programului. Sunt ilustrate tehnici de predicție a salturilor condiționate de la structuri adaptive corelate pe două niveluri la predictoare markoviene și neuronale precum și soluții de precalculare a adresei de salt. Tehnici aparte de predicție sunt aplicate categoriei de salturi și apeluri indirecte, cele mai de succes fiind predictoarele hibride, cascade sau metapredictoarele. Calitatea unui model de predicție depinde însă în mare măsură de calitatea datelor disponibile dar și de alegerea caracteristicilor schemei de predicție. Marea majoritate a schemelor se bazează pe utilizarea unui număr mai mare de caracteristici de intrare (cum ar fi adresa instrucțiunii de salt, comportamentul anterior din perspectivă globală a tuturor salturilor din program sau locală strict din perspectiva saltului supus predicției etc.) fără a lua în considerare cauza reală – salturile nepolarizate care sunt dificil de prezis, apar cu un comportament aproape aleator în același context de intrare, au o acuratețe de predicție scăzută și implicit generează pierderi de performanță. Bazat pe o abordare hibridă dintre Matematică și Calculatoare, am dezvoltat patru metrici referitoare la gradul de aleatorism al instrucțiunilor de salt nepolarizate: complexitatea Kolmogorov a programelor de calculator, rata de compresie, entropia discretă și acuratețea predicției folosind predictoare Markov cu straturi ascunse (Hidden Markov Models). Toate aceste metrici pot ajuta proiectantul de sisteme de calcul să înțeleagă mai bine predictibilitatea ramificațiilor de program și dacă predictorul poate fi îmbunătățit astfel încât să prezică cu o mai mare precizie salturile nepolarizate.

Capitolul al cincilea prezintă limitarea căii critice de program intitulată *issue bottleneck / data flow bottleneck* (limitarea consumatorului) împreună cu două tehnici anticipative și speculative propuse pentru depășirea acestei limitări și anume, reutilizarea dinamică a instrucțiunilor și respectiv

predicția dinamică a valorilor instrucțiunilor. Ambele soluții se bazează pe principiul vecinătății temporale a valorilor instrucțiunilor (*value locality*), formulat în 1996 de patru colective internaționale de cercetare de la AMD Nexgen, Carnegie Mellon, Wisconsin (USA) și Technion (Israel). Autorul acestei teze extinde conceptul asupra regiștrilor procesorului MIPS cu avantaje evidente, în primul rând reducerea drastică a numărului de predictoare implementate, dezvoltarea de predictoare hibride de valori (contextuale și incrementale) în care selecția predictorului curent se face prin intermediul unor selectoare dinamice de tip automat finit respectiv de tip neuronal, precum și ilustrarea unor considerații interesante asupra performanței globale de procesare. O idee de cercetare originală a constat în implementarea în mod selectiv a tehnicilor anticipative și speculative doar pentru valorile produse de instrucțiunile cu latență ridicată (Înmulțire / Împărțire sau instrucțiuni Load critice – care accesează cu miss primul nivel din cache și ajung în fruntea buffer-ului de reordonare), prezentându-se totodată și impactul la nivel de performanță de procesare și energie consumată.

În capitolul 6 sunt explorate oportunități de cercetare în arhitecturile multicore apărute ca o soluție la limitarea de tip *Power Wall* arătându-se că fertilizarea încrucișată între domeniul arhitecturilor de procesoare, calculul evoluționist cu metode de optimizare multi-obiectiv și folosirea regulilor fuzzy pentru reprezentarea cunoștințelor duce la un design eficient al procesorului. Obiectivul principal îl reprezintă optimizarea din perspectivă multiplă a arhitecturilor de procesare considerând ca metrici performanța, energia consumată, aria de integrare și temperatura disipată. Se motivează necesitatea implementării unor instrumente de explorare automată a spațiului de proiectare (*automatic design space exploration*) aferent procesoarelor mono/multicore cu peste 20 de parametri arhitecturali bazate pe algoritmi genetici de tip *Pareto* (*Non-dominated Sorting Genetic Algorithm* versiunea a 2-a), care după rularea unui număr consistent de generații permite analiza comparativă a rezultatelor folosind metrici de tipul: *front Pareto*, *hypervolum*, *coverage*, etc. O abordare originală o reprezintă îmbunătățirea algoritmului de explorare automată a spațiului de proiectare folosind cunoștințe de domeniu implementate sub forma unor reguli *fuzzy* în operatorii genetici care au sporit diversitatea și calitatea soluțiilor precum și viteza de convergență a procesului de explorare. Provocările actuale urmăresc îmbunătățirea procesului de explorare prin folosirea de tehnici de tip *Response Surface Modeling* (RSM) și identificarea modului de reprezentare adecvată a cunoștințelor de domeniu în sistemele multicore prin răspunsul la întrebările: care sunt regulile (mai importante), sunt acestea contradictorii, cum pot fi implementate în operatorii genetici de selecție, mutație sau încrucișare?

Capitolul 7 sintetizează planurile de viitor pe atât din perspectivă didactică, educațională cât și din perspectiva cercetării științifice, a creșterii capitalului uman din universitate și poate chiar cea a transferului tehnologic spre companii.