

Procesare Paralelă

Capitolul 1: Arhitecturi paralele

Bogdan Dumitrescu

Facultatea de Automatică și Calculatoare
Universitatea Politehnica București

Introducere

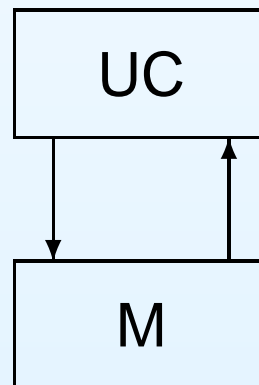
- Probleme de dimensiuni mari apar în multe domenii, atât științifice cât și tehnice
- Relativ puține tipuri de probleme; tipic: ecuații cu derivate parțiale
- Algebra liniară (sisteme de ecuații liniare, calculul valorilor proprii) apare frecvent
- Puterea totală de calcul disponibilă e uriașă, dar fiecare procesor e "relativ" lent; calculul paralel e soluția naturală
- Eficiența implementării este crucială, în special prin adecvarea algoritmilor la arhitecturile de calcul
- Arhitecturile evoluează rapid: cum adaptăm algoritmi ?

Cuprins

- Arhitecturi standard și caracteristicile lor
 - structură
 - operații favorizate
- Arhitecturi paralele
 - SIMD
 - MIMD (memorie comună, memorie distribuită)
 - Rețele de comunicație—topologii
- Cum adaptăm algoritmi la arhitecturi—strategii

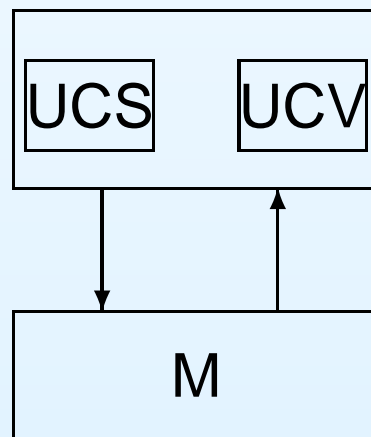
Arhitectura von Neumann

- Pentru fiecare operație, operanzii sunt aduși din memoria M în unitatea centrală UC, rezultatul fiind depus din nou în M
- Instrucțiunile stau și ele în M (alternativă: arhitectura Harvard, cu memorie separată pentru program)
- Unitatea de măsură standard pentru evaluarea timpului de execuție: operația în virgulă mobilă (flop)
- Model relativ corect până în anii '80



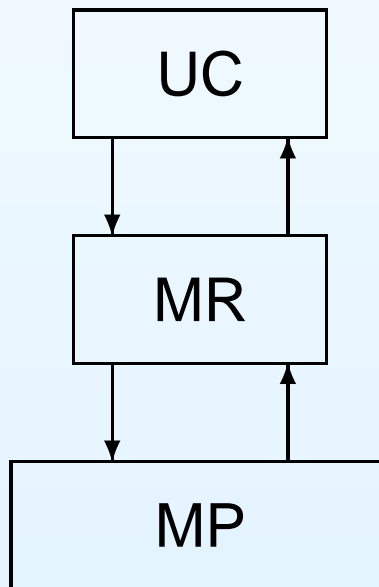
Arhitecturi vectoriale

- UCS—unitate de calcul obișnuită, pentru calcule scalare
- UCV—unitate de calcul vectorial, pe care operațiile vectoriale de genul $x + y$, $x, y \in \mathbb{R}^n$, se execută rapid
- UCV funcționează în regim pipeline, folosind și moduri favorizate de transfer din memorie (e.g. memoria e partajată în blocuri și există căi separate de alimentare a UCV pentru fiecare bloc)
- Operații eficiente: cu vectori de lungime suficient de mare



Arhitecturi cu memorie ierarhică

- Mai multe niveluri de memorie (cel puțin două)
- MR—memorie rapidă (cache), care poate alimenta UC cu date astfel încât să se efectueze un flop pe tact
- MP—memoria principală, în general mult mai lentă
- Operații eficiente: cele care refolosesc datele din MR



Memorie ierarhică—evoluție

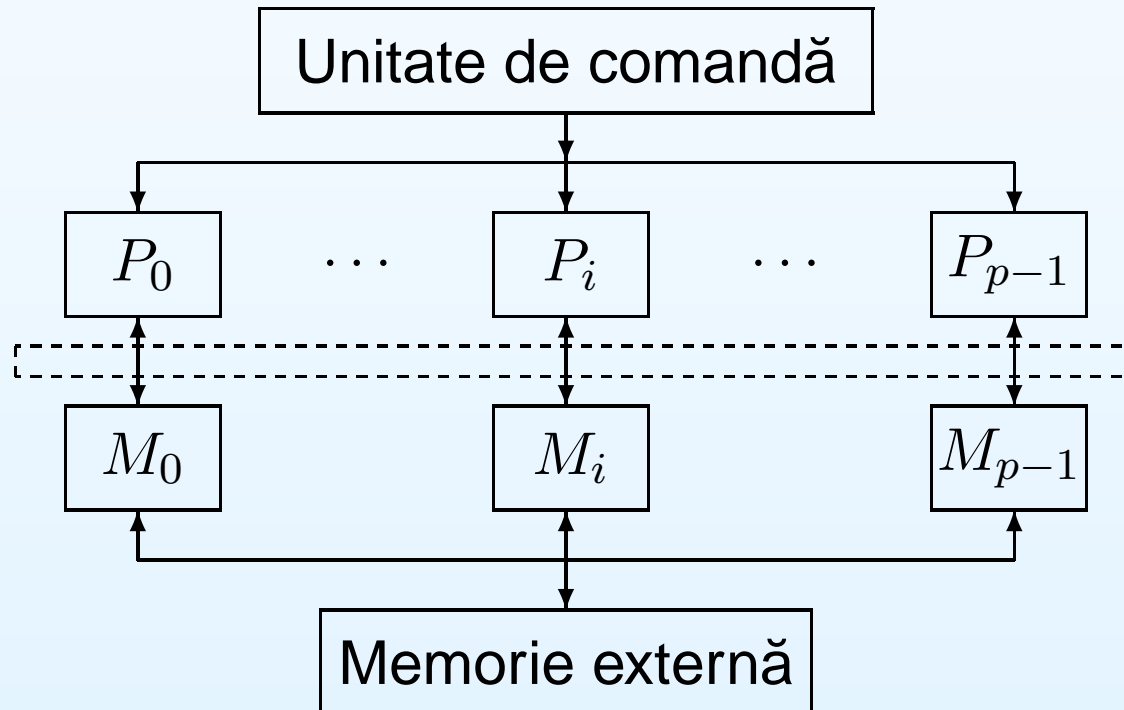
- Procesoare RISC (set redus de instrucțiuni), astfel încât decodificarea să fie rapidă
- Predicție a datelor (pe mai multe căi) și instrucțiunilor folosite, în încercarea de a aduce datele în MR înainte de a fi efectiv necesare
- Transferuri de blocuri de date din MP în MR
- Multicore—mai multe procesoare/UC, alimentate simultan (paralelism)
- Cheia eficienței rămâne refolosirea datelor (sau refolosirea imediată a rezultatelor), ceea ce necesită reorganizarea algoritmilor (operații la nivel de bloc)

Calculatoare paralele

- Tipuri principale: SIMD, MIMD
- SIMD (Single Instruction Multiple Data): aceeași instrucțiune se execută simultan asupra mai multor date
- MIMD (Multiple Instruction Multiple Data): mai multe fluxuri de instrucțiuni se execută asupra mai multor date (fiecare flux cu datele lui)
- Tipuri de MIMD
 - cu memorie comună
 - cu memorie distribuită

Arhitectura SIMD (1)

- Procesoarele P_i execută aceeași instrucțiune, dictată de unitatea de comandă
- Fiecare procesor își ia datele dintr-o memorie locală M_i

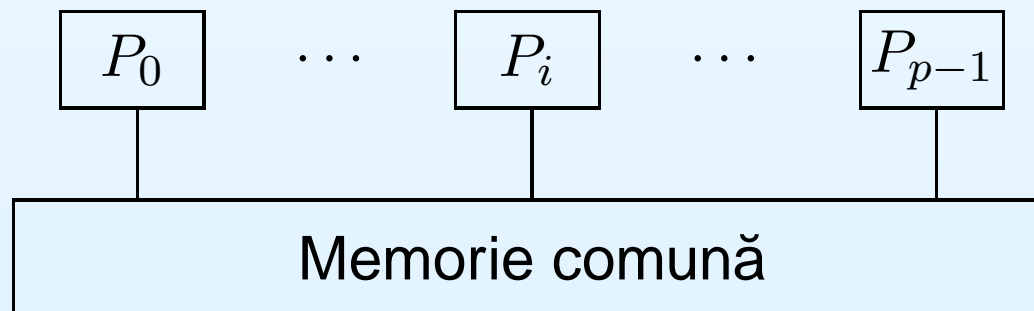


Arhitectura SIMD (2)

- Memoriile M_i sunt de obicei rapide și comunică cu o memorie principală (eventual externă)
- Între procesoarele P_i și memoriile M_i există o rețea de interconectare, care poate configura accesul (de exemplu: procesorul P_i ia datele din memoria M_{i+1})
- Procesoarele sunt de obicei foarte simple
- Operații favorizate: vectoriale, cu lungimea vectorului multiplu al numărului de procesoare
- Adunarea a doi vectori de lungime p se face într-un tact
- Problema critică: alimentarea cu date a $M_i \Rightarrow$ complicații arhitecturale

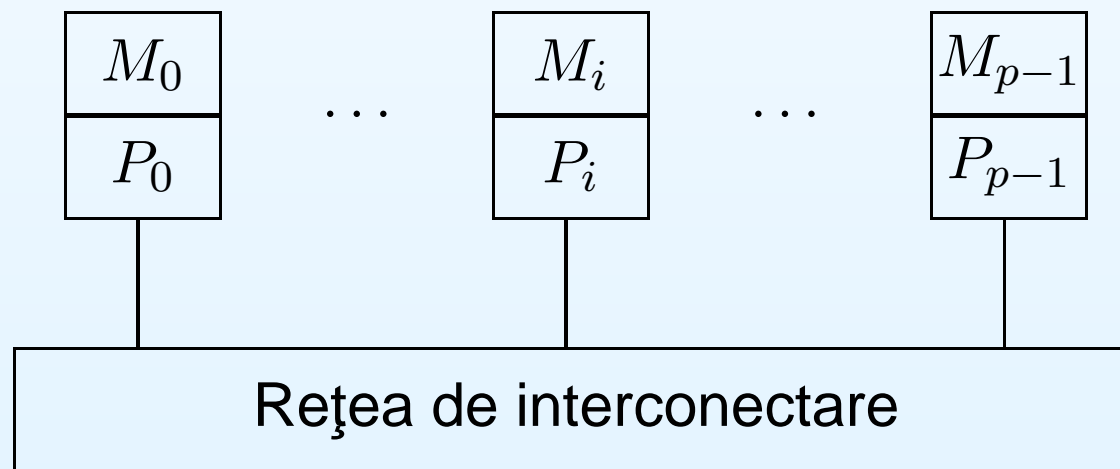
MIMD cu memorie comună

- Fiecare procesor execută instrucțiuni proprii
- Datele se află în memoria comună
- De obicei fiecare procesor are o memorie rapidă mică proprie (nereprezentată)
- Operații favorizate: paralele, la nivel de bloc
- Avantaj: comunicație simplă, prin intermediul memoriei
- Probleme: cu cât crește numărul de procesoare, cu atât crește probabilitatea conflictele de acces la memorie \Rightarrow scade viteza de calcul



MIMD cu memorie distribuită (1)

- Fiecare procesor are memorie proprie (arhitectura locală cu RISC și memorie ierarhică, de obicei)
- Comunicația se face printr-o rețea de comunicație, prin mesaje explicite
- Operații favorizate: paralele, la nivel de bloc
- Comunicația prin mesaje necesită algoritmi dedicați



MIMD cu memorie distribuită (2)

- Topologii ale rețelei de comunicație
 - fixă: inel, grilă (tor), hipercub, fat tree
 - configurabilă: switch, magistrale
- În calculatoarele actuale, topologia este de obicei transparentă pentru utilizator
- Biblioteci de funcții de comunicație
- Numărul de procesoare poate fi foarte mare, de ordinul miilor și peste (sute de mii)
- Putere de calcul maximă: > 10 petaflop/secundă
- Exercițiu: citiți www.top500.org

Grafuri

- O rețea de comunicație se modelează natural printr-un graf neorientat
- Procesoarele sunt nodurile, iar liniile de comunicație arcele
- Un *graf* neorientat este $G = (V, E)$, unde V e mulțimea nodurilor, iar $E \subset V \times V$ cea a arcelor
- Numărul de noduri ale unui graf: *ordin* (p , numărul de procesoare)
- Graf *simplu*: nu are bucle

Grafuri

- Două noduri sunt *adiacente* (vecine) dacă există un arc între ele
- *Gradul* unui nod: numărul de vecini ai nodului, adică numărul de arce plecând din nod
- *Gradul* unui graf (Δ): maximul gradelor nodurilor componente
- Graf *regulat*: gradele tuturor nodurilor sunt egale
- Numărul de arce al unui graf regulat este $p\Delta/2$

Grafuri

- *Drum* (sau *cale*) între două noduri $x, y \in V$: o secvență de noduri x_0, x_1, \dots, x_k , astfel încât două noduri consecutive sunt adiacente și $x_0 = x, x_k = y$.
- *Drum elementar*: nodurile componente x_1, \dots, x_{k-1} sunt distincte
- Într-o arhitectură cu memorie distribuită, un drum este o cale de comunicație între două procesoare
- *Lungimea* unui drum: numărul de arce ale sale, k în notația de mai sus
- *Distanța* dintre două noduri, $dist(x, y)$, este lungimea celui mai scurt drum între acestea
- *Diametrul* grafului (D): maximul tuturor distanțelor între perechi de noduri ale grafului, adică distanța dintre cele mai depărtate noduri.

Grafuri

- *Ciclu*: un drum elementar pentru care nodul inițial și cel final sunt identice
- *Ciclu hamiltonian*: un ciclu trecând (exact o dată) prin fiecare nod al grafului; lungimea sa este egală cu ordinul grafului
- Graf *conex*: există (cel puțin) un drum între oricare două noduri ale sale
- *Arbore*: graf conex fără cicluri

Deziderate topologice: hardware

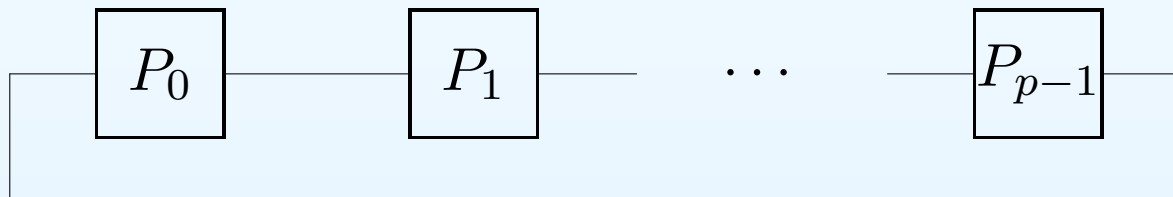
- Graf conex
- Grad mic și graf regulat: procesoare identice cu doar câteva canale de comunicație, deci ușor de realizat
- Număr total de arce relativ mic, fiecare arc însemnând o conexiune fizică realizată printr-un traseu electric pe o placă și, eventual, între plăci.

Deziderate topologice: software

- Diametru mic, astfel încât distanțele între procesoare să fie limitate, iar comunicarea cât mai facilă
- Cele două tipuri de deziderate sunt antagonice: cu cât gradul e mai mic, cu atât sunt mai puține arce, deci scad șansele să existe drumuri scurte între oricare două procesoare.
- Extremă: graful *total interconectat* (sau *complet*), în care există arce între oricare două noduri; $D = 1$, $\Delta = p - 1$, numărul de arce $p(p - 1)/2$
- La cealaltă extremă se află inelul

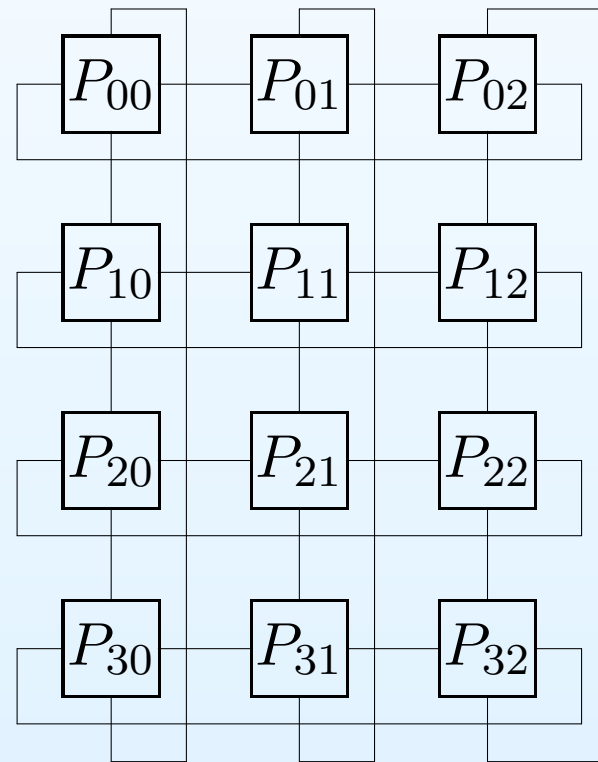
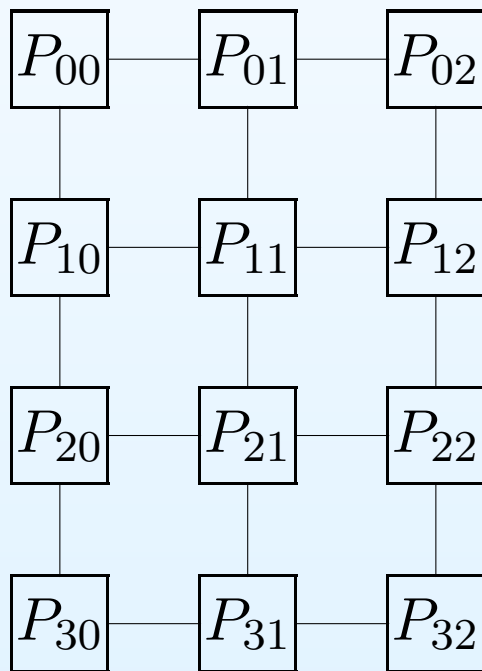
Inelul

- Procesorul P_i este legat de P_j , unde $j = (i - 1) \bmod p$ (P_j este vecinul din stânga) sau $j = (i + 1) \bmod p$ (P_j este vecinul din dreapta)
- Grad: $\Delta = 2$
- Numărul de arce: p
- Diametrul: $D = \lfloor p/2 \rfloor$



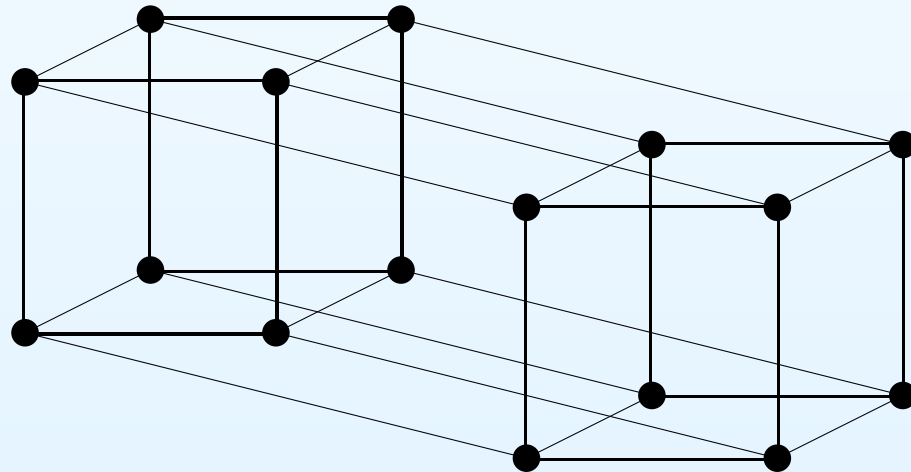
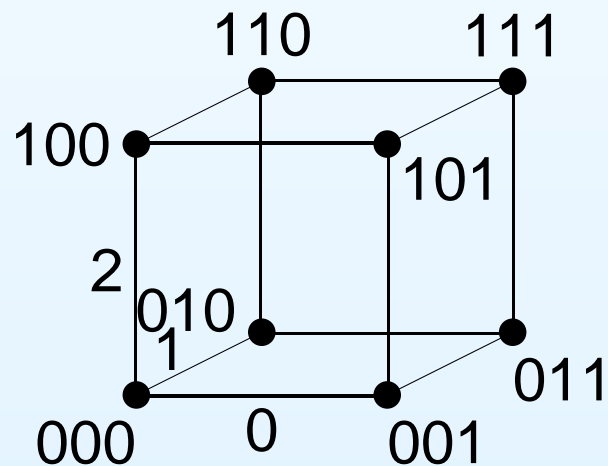
Grila și torul

- Pentru un tor pătrat $\sqrt{p} \times \sqrt{p}$
- Grad: $\Delta = 4$
- Numărul de arce: $2p$
- Diametrul: $D = 2\lfloor p/2 \rfloor$



Hipercubul

- Arhitectură foarte avantajoasă pentru unii algoritmi
- Notăm \mathcal{H}_d hipercubul de dimensiune d , care are $p = 2^d$ procesoare
- Fiecare procesor este legat de d vecini; P_i este conectat de P_j dacă reprezentările binare ale numerelor i și j diferă printr-un singur bit



Hipercubul

- Grad: $\Delta = d = \log p$
- Numărul de arce: $2pd = 2p \log p$
- Diametrul: $D = p$
- Distanța dintre două noduri este distanța Hamming dintre adresele lor
- Definiție recursivă: \mathcal{H}_d poate fi obținut luând două copii ale lui \mathcal{H}_{d-1} și adăugându-le arcele care unesc nodurile aflate în aceeași poziție în cele două copii; la adresele nodurilor se va adăuga un bit în poziția cea mai semnificativă: 0 pentru prima copie a lui \mathcal{H}_{d-1} , 1 pentru a doua

Recapitulare

Topologie	Δ	D	N
Inel	2	$O(p)$	$O(p)$
Grilă	4	$O(\sqrt{p})$	$O(p)$
Tor	4	$O(\sqrt{p})$	$O(p)$
Tor 3D	6	$O(\sqrt[3]{p})$	$O(p)$
Hipercub	$\log p$	$O(\log p)$	$O(p \log p)$

Cum obținem programe eficiente ?

- Scopul obișnuit în rezolvarea unei probleme de calcul științific este scrierea unui program cu timp de execuție cât mai mic
- (Avem în vedere probleme care se rezolvă de multe ori, cu date diferite de fiecare dată)
- Totuși, efortul de programare trebuie să fie rezonabil
- Dificultăți
 - un cod eficient pe un calculator poate fi foarte lent pe altul
 - arhitecturile sunt din ce în ce mai complexe, deci greu de modelat
 - arhitecturile evoluează repede
- Soluții: combinații ale celor expuse în paginile următoare

Soluția 1: compilatorul perfect

- Idealul programatorului:
 - se ia un algoritm optim dintr-un punct de vedere general (e.g. număr minim de operații)
 - se implementează într-un limbaj de nivel înalt
 - prin compilare se obține un cod (aproape) optim
- Compilatorul este optimizat pentru un anumit calculator
- Deși s-au făcut progrese, dificultăți majore:
 - descrierea precisă a arhitecturii
 - "înțelegerea" completă a programului de către compilator
 - timpul de compilare limitat (se pot obține coduri aproape optime pentru unele programe scurte, dar complexitatea compilării explodează pentru coduri mai lungi)
 - timpul de optimizare a unui compilator ar putea depăși durata de viață a calculatorului

Soluția 2: nuclee optimizate

- Se identifică operațiile "elementare" esențiale pentru un domeniu
- De exemplu, în algebra liniară, acestea sunt înmulțirea de matrice, rezolvarea de sisteme triunghiulare, plus altele similare cu complexitate mai mică (produs matrice-vector, produs scalar, etc.)
- Pentru fiecare calculator, se scriu biblioteci optimizate pentru aceste operații, e.g. BLAS
- Programele sunt scrise cu cât mai multe apeluri la biblioteci, astfel încât să fie portabile (după recompilare)
- Soluția cea mai răspândită acum, deși re-scrierea unei biblioteci optimale necesită multe luni de muncă pentru specialiști super-calificați
- Schimbări aparent mici în arhitectură pot necesita reoptimizarea bibliotecii

Soluția 3: adaptare automată

- Se păstrează ideea de bibliotecă cu funcții elementare
- Optimizarea bibliotecii pentru o nouă arhitectură se face automat, empiric
- Optimizarea se face prin intermediul unor parametri (e.g. dimensiuni de blocuri) într-unul sau mai multe coduri fixe sau prin generare de cod
- Pe baza unei proceduri de căutare, se măsoară timpii de execuție pentru diverse valori ale parametrilor
- Avantaj: optimizarea se face exact pe calculatorul țintă
- Dezavantaj: lipsă de robustețe la variații mari ale arhitecturii
- Exemplu: ATLAS (Automatically Tuned Linear Algebra Software)