

4. PLASAREA MODULELOR CU OBIECTIVUL ASIGURĂRII RUTABILITĂȚII CIRCUITELOR

4.1 Introducere

Plasarea este procesul de aranjare a componentelor unui circuit pe o suprafață. În cazul circuitelor FPGA, plasarea semnifică asignarea funcțiilor logice diferitelor celule ale circuitului, a căror poziție este fixă. Plasarea este o etapă importantă a procesului de proiectare, deoarece în această etapă se iau cele mai importante decizii. De exemplu, plasarea celulelor standard ale circuitelor VLSI are un impact major asupra vitezei și costului final al circuitului.

În proiectarea circuitelor, există diferite funcții obiectiv care trebuie minimizate, ca de exemplu lungimea conexiunilor, suprafața ocupată, puterea disipată, sau numărul conexiunilor care se intersectează. Practic, este imposibilă minimizarea simultană a tuturor acestor obiective, fiind necesară definirea unui set de criterii care reflectă aceste obiective. De obicei, este importantă minimizarea lungimii totale de rutare. Totuși, există și alte obiective care pot fi critice. De exemplu, în etapa de rutare, dacă problema este rutabilitatea circuitului, obiectivul poate fi modificat pentru a reflecta congestia conexiunilor în locul lungimii totale a acestora.

Atunci când este măsurată distanța între două module, există două metrici care pot fi utilizate: distanța Manhattan sau cea euclidiană. De exemplu, dacă pozițiile a două module pe o placă sunt (x_1, y_1) , respectiv (x_2, y_2) , distanța Manhattan este definită prin:

$$d = |x_1 - x_2| + |y_1 - y_2| \quad (4.1)$$

iar distanța euclidiană este definită prin:

$$d = (x_1 - x_2)^2 + (y_1 - y_2)^2 \quad (4.2)$$

Ca un exemplu, se consideră circuitul din Figura 4.1(a). Se presupune că este necesară plasarea porților pe o suprafață bidimensională. În Figura 4.1(b) este ilustrată o asemenea plasare, iar în Figura 4.1(c) este ilustrată aceeași plasare într-o formă simbolică. În cazul plasării simbolice, detaliile de rutare sunt omise. Din această plasare simbolică se poate obține însă o estimare a cerințelor de rutare. Considerând Figura 4.1(c), presupunem că rutarea unei conexiuni de la un modul la altul necesită o lungime a conexiunilor egală cu distanța Manhattan între module. De exemplu, lungimea conexiunii (3,6) este de 2 unități, lungimea totală a conexiunilor pentru plasarea din Figura 4.1(c) fiind de 10 unități. În Figura 4.1(d) se prezintă o altă plasare simbolică pentru care lungimea totală a conexiunilor ω este de 12. Figura 4.1(e) ilustrează o plasare unidimensională a circuitului, care necesită de asemenea o lungime a conexiunilor de 10 unități.

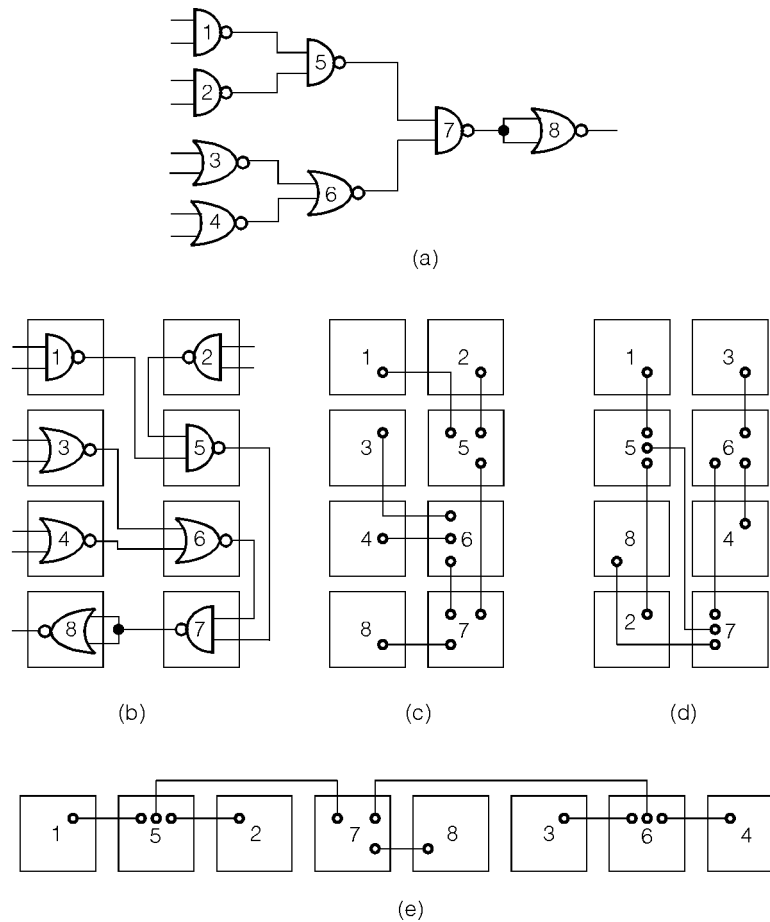


Figura 4.1. (a) Exemplu de circuit. (b) O plasare 2-D a porțiilor. (c) O plasare 2-D simbolică. (d) O plasare 2-D necesitând interconexiuni de 12 unități. (e) O plasare 1-D necesitând interconexiuni de 10 unități.

Lungimea totală a conexiunilor ω reprezintă o metrică utilizată pe scară largă pentru aprecierea calității plasării [146]. Considerăm plasarea simbolică din Figura 4.2(a). Același circuit poate fi plasat și în modul indicat în Figura 4.2(b), în ultimul caz suprafața necesară fiind însă mai mare. Această suprafață constă din două părți, cea funcțională și cea a conexiunilor. Suma suprafețelor celulelor funcționale reprezintă suprafața funcțională.

Pentru ambele plasări din Figura 4.2, suprafața funcțională este aceeași, fiind diferită numai suprafața conexiunilor. Aceasta din cauza unei distanțe minime care trebuie menținută între două conexiuni și între o celulă funcțională și o conexiune. De exemplu, considerăm conexiunile (2, 5) și (8, 9) din Figura 4.2(b). Aceste conexiuni trebuie plasate pe două piste verticale separate, între care trebuie menținută o distanță minimă.

O plasare care necesită un spațiu mare pentru conexiuni va necesita conexiuni de lungime mare, și deci lungimea totală a conexiunilor va avea o valoare mare. Deci, lungimea totală a conexiunilor ω este o metrică potrivită pentru suprafața ocupată de circuit. Avantajul utilizării acestei metrici este că se poate calcula simplu.

Plasarea modulelor cu scopul de a minimiza lungimea totală a conexiunilor ω este o problemă NP-completă [146]. Chiar și cazul cel mai simplu al problemei, și anume plasarea unidimensională, este dificil de soluționat. Pentru n module există un

număr de $\frac{n!}{2}$ plasări liniare. În practică, numărul modulelor sau al celulelor care trebuie plasate este foarte mare. De aceea, este impractică enumerarea tuturor configurațiilor de plasare posibile și selectarea celei mai bune. Au fost dezvoltate diferite tehnici euristice pentru rezolvarea problemei de plasare. Acestea nu permit de obicei obținerea soluției optime, dar cerințele de timp ale algoritmilor euristici sunt relativ reduse, reprezentând o funcție polinomială de n .

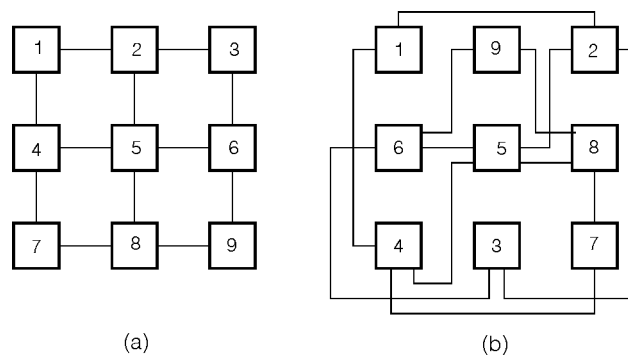


Figura 4.2. (a) Plasare optimă cu $\omega = 12$. (b) Soluție alternativă cu $\omega = 22$.

4.2 Definirea problemei de plasare

Fiind dată o colecție de module sau celule care dispun de porturi (pini de intrare, de ieșire, de alimentare și masă), dimensiunile acestor celule, și o colecție de conexiuni (reprezentând seturi de porturi care trebuie conectate), procesul de *plasare* constă din determinarea locațiilor fizice corespunzătoare pentru fiecare celulă. Aceste locații minimizează anumite funcții obiectiv, cu condiția respectării unor restricții impuse de proiectant, de procesul de implementare sau de stilul de proiectare. Exemple de restricții sunt evitarea suprapunerii celulelor sau cerința ca celulele să fie plasate într-o anumită suprafață rectangulară.

Într-un mod mai formal, problema de plasare poate fi definită astfel. Fiind dat un set de module $M = \{m_1, m_2, \dots, m_n\}$ și un set de semnale $S = \{s_1, s_2, \dots, s_k\}$, se asociază cu fiecare modul $m_i \in M$ un set de semnale S_{m_i} , unde $S_{m_i} \subseteq S$. Similar, cu fiecare semnal $s_i \in S$ se asociază un set de module M_{s_i} , unde $M_{s_i} = \{m_j \mid s_i \in S_{m_j}\}$. M_{s_i} se numește o conexiune de semnal. Este dat de asemenea și un set de locații $L = \{L_1, L_2, \dots, L_p\}$, unde $p \geq n$. Problema de plasare este de a asigna fiecare modul $m_i \in M$ unei locații unice L_j astfel încât să fie optimizat un anumit obiectiv.

În mod normal, fiecare modul este considerat zero-dimensional, și dacă m_i este asignat la locația L_j atunci poziția sa este definită prin coordonatele (x_j, y_j) . Uneori un subset al modulelor din M este fixat, deci asignat în prealabil unor locații, și numai modulele rămase pot fi asignate la locațiile rămase neasignate.

4.3 Funcții de cost și restricții

În cadrul procesului de proiectare, etapa de *plasare* este urmată de cea de *rutare*. O plasare este acceptabilă dacă se poate obține o rutare de 100% în cadrul suprafeței date. Funcția obiectiv care trebuie minimizată se poate scrie ca o sumă dintre

γ_1 și γ_2 . În cele mai multe cazuri, γ_1 este lungimea totală estimată a conexiunilor. În general, γ_2 reprezintă penalizări pentru soluțiile non-fezabile, fiind costul violării restricțiilor.

Executarea efectivă a rutării pentru a compara diferitele soluții de plasare nu este practică. De aceea, se utilizează diferite *estimări*. În secțiunile următoare se prezintă unele tehnici utilizate pentru estimarea lungimii conexiunilor necesare pentru o anumită plasare.

4.3.1 Estimarea lungimii conexiunilor

Viteza estimării are o influență fundamentală asupra performanțelor algoritmului de plasare. De aceea, procedura de estimare trebuie să fie cât mai rapidă. În plus, eroarea de estimare trebuie să fie aceeași pentru toate conexiunile.

O presupunere realistă pentru estimarea lungimii totale a conexiunilor este că rutarea utilizează geometria Manhattan [146]. Pentru o conexiune cu doi pini între modulul i și modulul j , lungimea Manhattan este $r_{ij} + c_{ij}$, unde r_{ij} și c_{ij} reprezintă numărul de linii și coloane care separă locațiile celor două module. Totuși, nu toate conexiunile sunt cu doi pini. Este necesară deci o metodă pentru a estima lungimea unei conexiuni multipunct. Există diferite tehnici utilizate, fiecare din acestea având avantaje și dezavantaje.

Metoda semi-perimetrului. Aceasta este o aproximare eficientă și foarte utilizată pentru a estima lungimea unei conexiuni. Metoda constă în determinarea celui mai mic dreptunghi care cuprinde toți pini conexiunii respective. Lungimea estimată a interconexiunilor este jumătatea perimetrului acestui dreptunghi. Pentru conexiuni cu doi și trei pini aceasta este o aproximare exactă. Această metodă asigură o estimare corespunzătoare pentru cea mai eficientă configurație de rutare, care este arborele Steiner. În cazul circuitelor foarte congestionate această metodă subestimează lungimea conexiunilor.

Metoda aproximării arborelui Steiner. Un *arbore Steiner* reprezintă calea cea mai scurtă pentru conectarea unui set de pini. În această metodă, pot exista ramificații din orice punct al unei conexiuni pentru conectarea la alți pini. Problema determinării arborelui Steiner minim este NP-completă. Există diferite metode pentru determinarea unei aproximări a arborelui Steiner [3] [45] [52] [84] [91] [99]. Se poate utiliza de asemenea algoritmul lui Lee, prezentat în Capitolul 5.

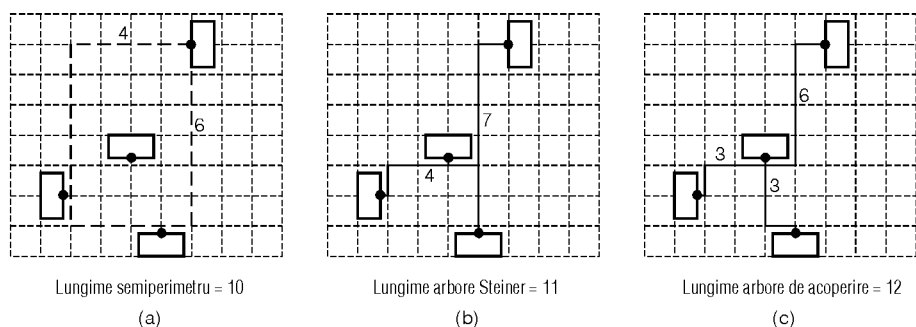


Figura 4.3. Aplicarea diferitelor metode pentru estimarea lungimii conexiunilor.

Metoda arborelui de acoperire minim. Spre deosebire de arborele Steiner, într-un arbore de acoperire minim ramificarea este permisă numai în pozițiile pinilor. Pentru o conexiune cu n pini, arborele poate fi construit prin determinarea distanțelor dintre toate perechile posibile de pini, și conectarea celor mai mici $(n - 1)$ muchii care

nu formează cicluri. Un algoritm cu o complexitate polinomială pentru determinarea arborelui de acoperire minim a fost elaborat de Kruskal.

Figura 4.3. prezintă exemple cu diferite metode de interconectare, pentru fiecare indicându-se lungimea conexiunilor.

4.3.2 Minimizarea lungimii totale a conexiunilor

Obiectivul principal al plasării este de a găsi o soluție care este complet rutabilă. De asemenea, suprafața ocupată de conexiunile de rutare trebuie să fie minimă. O metodă pentru a satisface această cerință este ca celulele puternic conectate să fie plasate aproape unele de altele. O funcție obiectiv a cărei minimizare se urmărește și care este des utilizată este $L(P)$, lungimea totală ponderată pentru toate conexiunile de semnal, exprimată ca:

$$L(P) = \sum_{n \in N} w_n \cdot d_n \quad (4.3)$$

unde:

d_n = lungimea estimată a conexiunii n ;
 w_n = ponderea conexiunii n ;
 N = setul de conexiuni.

În această estimare, lungimea fiecărei conexiuni este calculată în mod independent de celelalte conexiuni. De aceea, suprafața estimată este doar o aproximare grosieră a celei reale.

4.3.3 Minimizarea tăieturii maxime

Considerăm un spațiu pentru plasare sub forma unei suprafețe rectangulare (Figura 4.4), în care a fost plasat un circuit. Linia verticală de la $x = x_i$ divide suprafața într-o regiune din stânga L_i și o regiune din dreapta R_i . Față de această linie de tăietură, conexiunile se pot clasifica astfel:

- Conexiuni care se află în totalitate la stânga liniei de tăietură. Toți pinii unor asemenea conexiuni se vor afla în L_i .
- Conexiuni care se află în totalitate la dreapta liniei de tăietură. Toți pinii unor asemenea conexiuni se vor afla în R_i .
- Conexiuni care sunt tăiate de linie. Fiecare conexiune din această clasă va avea în mod necesar cel puțin un pin în L_i și cel puțin un pin în R_i .

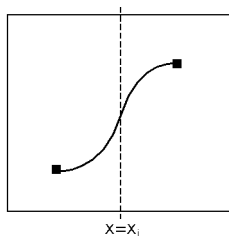


Figura 4.4. O linie de tăietură verticală și o conexiune intersectată de această linie.

Fie $\Phi_P(x_i)$ numărul de conexiuni de tipul c) pentru plasarea P tăiată de linia x_i . $\Phi_P(x_i)$ este o funcție de plasarea P . Pentru o anumită plasare P , fie $X(P)$ valoarea maximă a $\Phi_P(x_i)$ pentru fiecare i , deci:

$$X(P) = \max_i [\Phi_P(x_i)] \quad (4.4)$$

În mod similar, se pot defini linii de tăietură orizontale y_j și tăietura verticală maximă $Y(P)$ ca fiind:

$$Y(P) = \max_j [\Phi_P(y_j)] \quad (4.5)$$

Presupunem că plasarea se efectuează pentru o rețea de porți, și pentru o plasare dată, $X(P) = 10$ și $Y(P) = 15$. Aceasta înseamnă că, pentru o anumită linie de tăietură verticală v vor fi tăiate 10 conexiuni de această linie. În mod similar, pentru o anumită linie de tăietură orizontală h vor fi tăiate 15 conexiuni de această linie. Numărul de canale orizontale de-a lungul liniei de tăietură v trebuie să fie de cel puțin 10, deoarece în caz contrar circuitul nu va putea fi rutat utilizând plasarea P . În mod similar, trebuie să fie cel puțin 15 canale verticale de-a lungul liniei de tăietură h pentru a se putea realiza rutarea. Existența acestui număr de canale nu garantează faptul că rutarea va putea fi executată complet, fiind deci o condiție necesară, dar nu și suficientă. Din cele de sus rezultă că $X(P)$ și $Y(P)$ sunt strâns legate de rutabilitatea rețelei de porți. Dacă se dă o rețea de porți cu H_{max} canale orizontale și V_{max} canale verticale pe o linie de caroiaj, atunci trebuie găsită o plasare pentru care $X(P) \leq H_{max}$ și $Y(P) \leq V_{max}$.

Tăieturile $\Phi_P(x_i)$ și $\Phi_P(y_j)$ sunt de asemenea strâns legate de lungimea totală a conexiunilor $L(P)$. Presupunând că spațierea liniilor de caroiaj este de 1 unitate, se poate arăta că

$$L(P) = \sum_i \Phi_P(x_i) + \sum_j \Phi_P(y_j) \quad (4.6)$$

unde însumarea se efectuează pentru toate liniile de tăietură posibile.

Din cele expuse rezultă că reducerea tăieturii orizontale $X(P)$ și a celei verticale $Y(P)$ prin selectarea unei plasări P corespunzătoare poate mări probabilitatea rutării rețelei de porți. În plus, minimizarea $X(P)$ și $Y(P)$ poate avea de asemenea o influență benefică asupra lungimii totale a conexiunilor $L(P)$.

În cazul circuitelor care utilizează celule standard sau macro-celule, minimizarea $X(P)$ și $Y(P)$ este de asemenea importantă pentru a se reduce lungimea totală a conexiunilor.

4.3.4 Minimizarea densității maxime

O altă măsură a rutabilității unei plasări este *densitatea* $D(P)$, definită după cum urmează. Presupunem că suprafața disponibilă pentru plasare este divizată într-o grilă. Figura 4.5(a) ilustrează o rețea de porți cu trei linii și trei coloane. Blocurile funcționale sunt hașurate. În Figura 4.5(b) se arată separat o porțiune A a suprafeței rezervate pentru conexiuni (un bloc de interconectare). Prin această regiune poate trece un număr fix de conexiuni orizontale, care se numește capacitate orizontală a regiunii. În mod similar, se definește și o capacitate verticală a blocului de interconectare. Figura 4.5(b) prezintă de asemenea un canal B care permite conexiuni verticale. Pentru acest canal este definită o capacitate verticală.

Fiind dată o plasare P , este posibilă estimarea numărului de conexiuni care trebuie să treacă prin fiecare muchie e_i a unui canal (sau bloc de interconectare). Dacă această estimare este $\eta_P(e_i)$, iar $\psi_P(e_i)$ este capacitatea muchiei e_i , atunci se poate defini densitatea muchiei e_i ca fiind

$$d_P(e_i) = \frac{\eta_P(e_i)}{\psi_P(e_i)} \quad (4.7)$$

Această densitate $d_p(e_i)$ trebuie să fie mai mică decât 1 (sau cel mult 1) în scopul rutabilității. Măsura rutabilității plasării este dată de:

$$D(P) = \max_i [d_p(e_i)] \quad (4.8)$$

unde maximul se calculează pentru toate muchiile e_i ale regiunilor de rutare (blocuri de interconectare).

4.3.5 Maximizarea performanțelor

Pe măsura avansului tehnologic, viteza de comutare a tranzistoarelor este din ce în ce mai ridicată. Ca urmare a acestei tendințe, întârzierile datorate interconexiunilor devin importante comparativ cu cele datorate comutării tranzistoarelor. În cazul unor tehnologii ca ECL, acest efect este pronunțat.

Deoarece timpul de comutare al porților logice a scăzut la ordinul picosecundelor, frecvența ceasului circuitelor VLSI a devenit din ce în ce mai dependentă de propagarea semnalelor prin interconexiunile circuitelor. Pentru numeroase circuite complexe, întârzierile datorate interconexiunilor contribuie cu mai mult de jumătate la ciclul de ceas, și porțiunea datorată timpului de propagare din cadrul ciclului continuă să crească. Acest fapt are un impact major asupra succesului procesului de proiectare. În etapele de început ale procesului de proiectare nu este posibilă determinarea frecvenței ceasului numai pe baza caracteristicilor logice și a întârzierilor de comutare a circuitelor. Timpii mari de propagare, care se datorează caracteristicilor electrice ale interconexiunilor, nu permit obținerea frecvenței de ceas așteptate.

Pentru verificarea și îmbunătățirea proprietăților temporale ale circuitelor, proiectanții utilizează utilitare numite analizoare de întârziere sau verificatoare de timp. Acestea permit testarea problemelor care pot apărea pe căile scurte sau lungi ale semnalelor.

În timpul procesului de proiectare, întârzierile de propagare ale interconexiunilor nu sunt cunoscute înainte de amplasarea modulelor. Problemele de temporizare pe căile lungi care apar după amplasare sunt foarte dificil de corectat deoarece ele pot necesita nu numai noi iterații ale etapelor de proiectare fizică, ci uneori, un număr mare de iterații ale etapei de proiectare logică.

Au fost sugerate trei metode generale pentru corectarea problemelor de temporizare pe căile lungi [146]. Prima metodă constă în efectuarea unor modificări în cadrul logicii. De exemplu, întârzierea unei căi poate fi redusă în mod substanțial prin reducerea încărcării unor elemente de circuit de pe această cale.

A doua metodă constă în dimensionarea tranzistoarelor pentru a crește viteza unor elemente de circuit pe calea lentă. Prin creșterea dimensiunii unora din tranzistoare, întârzierile de comutare, ca și întârzierile de propagare prin conexiunile pe care se află aceste tranzistoare pot fi reduse substanțial.

O metodă de a crește viteza unui circuit fără a modifica proiectarea logică a acestuia constă în reducerea timpului de propagare la minim. Acest scop poate fi atins prin impunerea unor constrângeri de temporizare pentru interconexiunile și căile circuitului. A treia metodă utilizează această tehnică. Au fost raportate diferite încercări pentru ca proiectarea fizică să țină cont de cerințele de temporizare. Un circuit VLSI care este optimizat din punct de vedere al temporizării permite o creștere cu 25-35% a frecvenței ceasului, fără efectuarea unor modificări în proiectarea logică.

Performanțele de viteză ale unui circuit pot fi caracterizate prin întârzierea combinatională cea mai mare de la un pin de ieșire la un pin de intrare. Dacă întârzierile căilor semnalelor trebuie menținute sub o anumită valoare maximă, atunci trebuie verificate întârzierile conexiunilor. Deoarece plasarea afectează cerințele de interconectare ale circuitului, obiectivul unei probleme de plasare poate fi modificat pentru a satisface cerințele de temporizare ale căilor semnalelor.

4.4 Sinteza metodelor de plasare

După cum s-a arătat în secțiunea 4.1, problema de plasare este NP-completă. Funcțiile de cost și restricțiile prezentate în secțiunea precedentă nu îmbunătățesc această situație. Pentru probleme complexe, spațiul soluțiilor nu permite enumerarea tuturor soluțiilor, motiv pentru care se utilizează tehnici euristice care necesită timp redus și de execuție și pot găsi soluții apropiate de cele optime.

Metodele euristice pentru plasare se pot clasifica în două categorii: *constructive* și *iterative*.

Metodele constructive realizează plasarea câte unei celule la un moment dat, construind astfel soluția finală. Aceste metode sunt de două tipuri: metode bazate pe partiționare și metode analitice. În ultimul timp, au fost obținute soluții de calitate prin combinarea ambelor strategii [67]. Metodele constructive sunt de obicei rapide și generează soluții de calitate. Totuși, acestea sunt în general limitate în privința alegerii obiectivelor, și de multe ori nu obțin optimul global al problemei de plasare.

Metodele iterative au ca scop îmbunătățirea soluțiilor existente, în special a celor obținute prin metode constructive. În general, într-un pas iterativ acestea selectează o subproblemă locală de dimensiuni mici pentru a fi rezolvată printr-o metodă exactă sau euristică. Aceste metode se împart de asemenea în două categorii, după cum ele aplică tehnici deterministice sau aleatoare.

Metodele iterative deterministice oferă de obicei proiectantului alegerea obiectivelor de plasare, dar de multe ori aceste metode obțin doar un minim local. Timpul necesar obținerii soluției este mult mai redus decât în cazul metodelor aleatoare.

Metodele iterative aleatoare pot accepta, cu o anumită probabilitate redusă, și soluții de plasare intermediară de calitate inferioară. Astfel, acestea au posibilitatea de a evita soluțiile optime locale și de a se apropia de optimul global, dacă dispun de un timp de calcul suficient. Există două metode aleatoare de bază, evoluția simulată și călirea simulată.

În literatură a fost publicat un număr mare de algoritmi pentru plasare [37], [57] [67], [68], [71], [82], [87], [88], [104], [109], [113], [119], [124], [139], [141], [151], [152], [153], [189]. Metodele bazate pe partiționare implică aplicarea recursivă a unui algoritm de partiționare, de obicei algoritmul Kernighan-Lin sau algoritmul de călire simulată [102]. Pachetul de programe *TimberWolf* [151] pentru plasare și rutare, care se bazează pe metoda de călire simulată, a fost primul program care a utilizat această metodă pentru problema de plasare. Rezultatele care au fost obținute sunt foarte bune pentru un număr de câteva mii de celule. Deși metoda de călire simulată poate obține soluția optimă globală, timpul de calcul necesar pentru circuite de dimensiuni mari este foarte mare, și de multe ori nu este acceptabil în practică. Pentru a elimina dezavantajul complexității ridicate a algoritmului de călire simulată, au fost propuse mai multe tehnici de creștere a vitezei pentru acest algoritm [107], [119].

Lam și Delosme [107] au elaborat o schemă îmbunătățită de călire, care păstrează rata de acceptare a călirii simulate la o valoare constantă de 44% pe parcursul domeniului intermediar de temperatură. Cu această schemă au fost raportate îmbunătățiri atât ale performanțelor, cât și a timpului de calcul. Această îmbunătățire a fost inclusă în versiunile ulterioare ale programului *TimberWolfSC*.

Mallela și Grover [119] au propus o călire simulată în două etape, utilizând o strategie de grupare de jos în sus. În prima etapă, se aplică o tehnică de grupare. Celulele cu o proximitate mai mare în privința numărului de terminale și a numărului de conexiuni sunt combinate pentru a forma un grup. Ca rezultat al grupării, spațiul de căutare al călirii simulate este redus în principiu cu dimensiunea grupurilor. Mallela și Grover raportează o reducere cu un factor de 2 până la 3 a timpului de calcul, și o reducere cu 6-17% a lungimii conexiunilor, măsurată prin metoda semi-perimetrului.

Această tehnică de călire simulată bazată pe grupare a fost inclusă în sistemul de plasare LTX2 al Bell Laboratories.

Hamada *et al.* [88] au propus o tehnică de plasare pe baza partiționării ierarhice prin algoritmul de tăietură proporțională elaborat de Cheng și Wei [50]. Prin aplicarea recursivă a partiționării, circuitul este descompus într-un arbore de grupuri, iar problema inițială de plasare este soluționată printr-o secvență de procese de călire simulată, unde fiecare grup este tratat ca o super-componentă. Această metodă reduce în mod considerabil complexitatea problemei.

Pentru soluționarea problemei de plasare prin tehnici iterative a fost propusă aplicarea metodei de asignare liniară, o metodă exactă și eficientă din punct de vedere computațional [67]. Această metodă a fost utilizată pentru translatarea unei plasări globale conținând celule suprapuse într-o plasare finală prin minimizarea distanței cu care celulele sunt mutate din pozițiile lor suprapuse. Algoritmii de asignare liniară au fost propuși de asemenea pentru soluționarea problemei speciale de plasare în care toate celulele au aceeași dimensiune și sunt specificate locațiile posibile ale acestora.

Pentru rezolvarea problemei de plasare liniară, care este o problemă fundamentală pentru proiectarea circuitelor VLSI, au fost propuse diferite metode. Li *et al.* [109] au propus o metodă spectrală în care se utilizează o funcție obiectiv rezultată dintr-o combinație a unei funcții liniare cu o funcție quadratică. Se utilizează astfel atât avantajul unei funcții obiectiv liniare, care permite obținerea unei plasări de calitate mai bună în privința lungimii conexiunilor, cât și avantajul unei funcții cuadractice, care tinde să plaseze componentele mai dispersat, rezultând o soluție mai fezabilă.

Problema de plasare poate fi transformată într-o problemă de optimizare numerică. Hanan și Kurtzberg au redus problema de plasare la cea de rezolvare a unui set de ecuații liniare pentru a determina locațiile de echilibru (coordonatele ideale x , y) ale celulelor [146].

Plasarea orientată pe performanțe a fost studiată de numeroși autori, în special plasarea cu restricții de timp. Metodele raportate pot fi clasificate în trei categorii [146]. Prima categorie transformă restricțiile de timp de pe căile critice în ponderi ale conexiunilor. Aceste ponderi sunt utilizate pentru a categorisi conexiunile și a influența procedura de plasare. A doua categorie transformă restricțiile de timp ale căilor în limite de timp ale conexiunilor. Aceste limite de timp sunt convertite în limite de lungime ale conexiunilor și sunt furnizate programului de plasare, care încearcă satisfacerea acestora. A treia metodă constă în furnizarea pentru programul de plasare a unui set de căi critice, împreună cu cerințele lor de timp. Aceste căi sunt monitorizate în timpul procesului de plasare.

Problema de plasare a circuitelor FPGA a fost abordată prin diferite metode, dintre care și prin metoda călirii simulate, într-un mod similar cu plasarea celulelor standard. În timp ce tehnicile elaborate pentru celulele standard sunt suficiente pentru acele circuite FPGA la care o mare porțiune a spațiului de pe cip este dedicată resurselor de rutare [178], în cazul arhitecturilor FPGA cu resurse limitate de rutare sunt necesare tehnici speciale. Ebeling *et al.* [71] au elaborat un program de plasare pentru circuitele FPGA Triptych, program care se bazează pe metoda de călire simulată. Beetem [18] a propus un algoritm de îmbunătățire iterativă pentru plasarea și rutarea simultană a circuitelor FPGA. Nag și Roy [125] a prezentat un algoritm incremental de plasare pentru circuitele FPGA cu o arhitectură bazată pe rânduri de celule, care analizează informațiile de întârziere a semnalelor și de rutabilitate pentru a obține plasări de calitate mai bună. Togawa *et al.* [167] au propus o metodă pentru plasarea și rutarea simultană a circuitelor FPGA simetrice, metodă bazată pe bipartiționarea ierarhică.

Gao [82] a elaborat algoritmi de plasare orientați pe performanțe bazați pe conexiuni și pe căi de rutare, pentru rețele de porți, macro-celule și circuite FPGA Xilinx, cu scopul minimizării întârzierii semnalelor la pinii de ieșire și a nesimetriei semnalelor la intrările modulelor. În cazul algoritmului bazat pe conexiuni, cerințele

de întârziere sunt translatate mai întâi în constrângeri de proiectare fizică, cum sunt constrângeri ale conexiunilor. Algoritmul de plasare generează apoi o plasare ghidată de aceste constrângeri. În cazul algoritmului bazat pe căi de rutare, întârzierile căilor sunt considerate în mod explicit în timpul procesului de plasare. Acest algoritm încearcă să minimizeze lungimea totală a conexiunilor și timpii de întârziere la pinii de ieșire.

Pentru rezolvarea problemei de plasare au fost propuse diferite metode neconvenționale, ca rețele neuronale, algoritmi genetici, logica fuzzy, sau prelucrarea paralelă. Yu [187] a modificat rețelele neuronale introduse de Hopfield și Tank pentru rezolvarea problemei de plasare. Funcția de energie utilizată de Hopfield are mai multe minime, dintre care unele sunt minime locale; rețeaua neuronală poate converge în oricare din acestea. În plus, este dificilă determinarea parametrilor rețelei. Rezultatele obținute de Yu nu au fost satisfăcătoare. Unele din problemele apărute sunt timpii mari de simulare, calitatea redusă a soluțiilor și dependența soluțiilor de parametrii rețelei.

Algoritmi genetici au fost de asemenea utilizați pentru plasarea celulelor circuitelor VLSI [124], [152]. Mohan și Mazumder [124] au elaborat un algoritm genetic pentru plasarea celulelor standard. A fost implementată atât o versiune serială a algoritmului, cât și una paralelă, care rulează pe o rețea de stații de lucru. Creșterea de viteză obținută este liniară cu numărul de procesoare utilizate. Algoritmul paralel păstrează calitatea versiunii seriale.

Plasarea este o problemă cu obiective multiple, unde numeroase decizii care sunt luate în timpul căutării soluției sunt calitative. O abordare potrivită a acestei categorii de probleme este utilizarea logicii fuzzy. O descriere a problemei de plasare care se bazează pe logica fuzzy a fost publicată de Lin și Shragowitz [113].

Au fost realizate diferite implementări paralele ale metodei de călire simulată pe diferite tipuri de calculatoare paralele: cu memorie partajată, cu transmitere de mesaje și memorie locală, și calculatoare cu paralelism masiv, ca de exemplu Connection Machine. Kravitz și Rutenbar [104] au prezentat un algoritm de călire simulată pentru celule standard, care a fost implementat pe un multiprocesor. Aceștia au arătat că algoritmul de călire simulată poate fi accelerat în două moduri pe un multiprocesor cu memorie partajată: prin executarea mai multor mutări în paralel, și prin executarea în paralel a prelucrărilor necesare pentru fiecare mutare.

Rose *et al.* [141] au conceput euristici pentru plasarea paralelă a celulelor standard cu o calitate echivalentă cu cea a călirii simulate. Aceștia au utilizat o metodă rapidă bazată pe tăietura minimă pentru a evita partea de călire lentă la temperaturi înalte. În faza de călire la temperaturi joase, spațiul din cadrul circuitului este partiționat și este asignat diferitelor procesoare, astfel încât fiecare procesor mută celule într-o anumită zonă și, ori de câte ori o mutare este acceptată, transmite rezultatul tuturor procesoarelor.

În continuare se descriu într-un mod mai detaliat unele metode de plasare. În secțiunea 4.4.1 se prezintă plasarea prin metode constructive. În secțiunea 4.4.2 se descrie un algoritm de plasare care utilizează în mod repetat metoda de partiționare pe baza tăieturii minime. În secțiunea 4.4.3 se descrie modul de utilizare a algoritmului de călire simulată pentru plasare, și se prezintă algoritmul de plasare utilizat de pachetul de programe TimberWolf. În secțiunea 4.4.4 se prezintă plasarea prin partiționarea ierarhică, în care se utilizează metoda de partiționare pe baza tăieturii minime. În secțiunea 4.4.5 se prezintă plasarea prin metode de optimizare numerică, fiind descrisă o metodă numită plasare controlată de forțe. În secțiunea 4.4.6 se prezintă utilizarea metodelor spectrale pentru rezolvarea problemei de plasare liniară. În secțiunea 4.4.7 se descriu pe scurt alte metode iterative de plasare: asignarea quadratică, optimizarea rețelelor rezistive, algoritmul Steinberg, metoda spațiului grafurilor, și plasarea bazată pe operații cu linii și coloane.

4.4.1 Plasarea constructivă inițială

Un algoritm constructiv generează o configurație de plasare completă numai la sfârșitul întregului proces. Un asemenea algoritm se utilizează adesea pentru generarea unei plasări inițiale, urmând ca aceasta să fie îmbunătățită printr-o metodă iterativă. Plasarea inițială este esențială pentru obținerea unei soluții optime. Există diferite soluții optime, corespunzătoare diferitelor plasări inițiale.

Se consideră ca exemplu următorul algoritm pentru plasarea a n module într-o linie. Suprafața pe care se amplasează modulele este divizată în n locații, inițial fiecare locație fiind liberă. În fiecare iterație a algoritmului, se va plasa un modul într-una din locațiile libere. La sfârșitul unei iterații, va exista o plasare parțială a unui subset de module. Există două decizii care trebuie luate în fiecare iterație. Acestea se referă la selectarea modulului care va fi adăugat la plasarea parțială și la locația în care se va plasa modulul selectat.

Se pot utiliza diferite euristici pentru deciziile care trebuie luate. De exemplu, o euristică posibilă pentru selecție este alegerea celui mai puternic conectat cu plasarea parțială existentă. Presupunând că plasarea parțială este formată din modulele m_1, m_2, \dots, m_i , se examinează fiecare din modulele neplasate m_j și se calculează cantitatea

$$A_{mj} = \sum_{k=1}^i c_{m_j m_k} \quad (4.9)$$

unde $c_{m_j m_k}$ reprezintă conectivitatea între modulul neplasat m_j și un modul plasat m_k . Astfel, A_{mj} indică numărul de conexiuni de la m_j la modulele deja plasate $\{m_1, m_2, \dots, m_i\}$. Se va selecta modulul pentru care A_{mj} este maxim. Această strategie este cunoscută sub numele de strategie de *conectivitate maximă*.

Modulul selectat poate fi plasat în oricare din cele $(n - i)$ locații libere. Se poate estima modificarea funcției de cost pentru fiecare din cele $(n - i)$ opțiuni, alegând opțiunea care este cea mai avantajoasă. De exemplu, dacă funcția de cost este lungimea totală a conexiunilor, este selectată locația care va determina creșterea minimă a lungimii conexiunilor.

Pentru a obține o plasare parțială inițială, se poate selecta un singur modul, în mod aleator sau pe baza unui criteriu euristic. Acest modul poate fi plasat, de exemplu, în locația din mijloc. Este logic ca inițial să fie selectat modulul cu conectivitatea maximă.

În continuare se vor descrie într-un mod mai formal cele două etape ale unui algoritm de plasare constructivă.

Pentru selecția unui modul neplasat, fie X setul modulelor plasate, set în care unele module sunt fixe (module hard), ca de exemplu celule de I/E, și fie Y setul modulelor neplasate. În mod intuitiv, următorul modul candidat din setul Y trebuie să fie cel cu cele mai multe conexiuni cu modulele din setul X , pentru a se obține lungimea de rutare cea mai mică. Regulile de selecție a modulului candidat vor fi generate pe baza matricii de capacitate C , după cum urmează:

$$M_Y = \max_{y \in Y} \left\{ \sum_{x \in X} c_{yx} \right\} \quad (4.10)$$

unde x este un modul din setul X al modulelor plasate, iar y este un modul care trebuie plasat, aparținând setului Y al modulelor neplasate. În cazul unei indecizii, deci atunci când poate fi ales mai mult decât un candidat, se poate selecta în mod arbitrar oricare din acei candidați.

De exemplu, considerăm următoarea matrice de capacitate \mathbf{C} pentru modulele A, B, C, D și E :

$$\mathbf{C} = \begin{bmatrix} 0 & 4 & 2 & 3 & 5 \\ 4 & 0 & 0 & 5 & 4 \\ 2 & 0 & 0 & 0 & 1 \\ 3 & 5 & 0 & 0 & 2 \\ 5 & 4 & 1 & 2 & 0 \end{bmatrix} \quad (4.11)$$

Dacă modulele A, B și E au fost plasate, următorul modul care trebuie plasat este D , deoarece $M_D = 3 + 5 + 2 = 10 > M_C = 2 + 1 = 3$.

Mai general, dacă se consideră că ponderea este un factor important al setului de semnale, se poate introduce o pondere w_{ij} pentru conexiunea dintre modulele i și j . O formă generală a regulii de selecție este următoarea:

$$M_Y = \max_{y \in Y} \left\{ \sum_{x \in X} q_{yx} \right\} \quad (4.12)$$

unde $\mathbf{Q} = [q_{ij}]$, $q_{ij} = c_{ij} w_{ij}$. Dacă $w_{ij} = 1$, nu există deosebire între matricile \mathbf{C} și \mathbf{Q} .

Există diferite alte reguli de selecție [189]. De exemplu, o altă regulă simplă de selecție poate fi bazată pe matricea de capacitate, alegând ca și candidat modulul cu

$$\max_i \left\{ \sum_j c_{ij} \right\}, \quad \forall j \in V \quad (4.13)$$

unde modulul i aparține setului de module neplasate, iar V este setul de noduri ale grafului G al circuitului.

Pentru plasarea modulului selectat, trebuie găsite coordonatele unui punct optim pentru care lungimea totală de rutare este minimă. După cum s-a indicat în secțiunea 4.1, lungimea de rutare poate fi măsurată prin distanța euclidiană sau prin distanța Manhattan. Dacă există m module fixe M_j aflate în pozițiile (x_j, y_j) , unde x_j și y_j sunt întregi, $j = 1, 2, \dots, m$, și un modul M în poziția (x, y) care poate fi mutat, cu o capacitate c_{Mj} între M și modulul fix M_j , $j = 1, 2, \dots, m$, lungimea totală de rutare a modulului M la toate modulele fixate poate fi exprimată în metrica euclidiană astfel:

$$D = \sum_{j=1}^m c_{Mj} [(x - x_j)^2 + (y - y_j)^2] \quad (4.14)$$

iar în metrica Manhattan astfel:

$$D = \sum_{j=1}^m c_{Mj} [|x - x_j| + |y - y_j|] \quad (4.15)$$

Pentru a găsi punctul optim (x_{opt}, y_{opt}) pentru care lungimea totală de rutare D este minimă, se vor prezenta separat cazurile celor două metrici.

A. Metrica euclidiană. Se observă că în ecuația (4.14), cele două componente din membrul drept sunt independente. Deci, lungimea totală în direcția x este independentă de lungimea totală în direcția y , care este cea mai importantă pentru ordonarea liniară sau plasare. Astfel, cele două componente pot fi considerate separat. Se consideră lungimea totală în direcția x :

$$D_x = \sum_{j=1}^m c_{Mj} (x - x_j)^2 \quad (4.16)$$

Dacă se calculează derivata lui D_x față de x și se egalează cu zero, se obține:

$$x_{opt} = \frac{\sum_{j=1}^m c_{Mj} x_j}{\sum_{j=1}^m c_{Mj}} \quad (4.17)$$

care este punctul optim pentru obținerea valorii minime pentru D_x . Similar,

$$y_{opt} = \frac{\sum_{j=1}^m c_{Mj} y_j}{\sum_{j=1}^m c_{Mj}} \quad (4.18)$$

Se definește:

$$x_{mopt} = [x_{opt} + 1] \quad (4.19)$$

$$y_{mopt} = [y_{opt}] \quad (4.20)$$

dacă modulul M_m este în legătură numai cu unul din modulele plasate, iar în caz contrar:

$$x_{mopt} = [x_{opt} + 0.5] \quad (4.21)$$

$$y_{mopt} = [y_{opt} + 0.5] \quad (4.22)$$

unde (x_{mopt}, y_{mopt}) este numit *punct optim local* pentru modulul M_m , iar $[x]$ este cel mai mare întreg mai mic sau egal cu x .

B. Metrica Manhattan. Înainte de a determina punctul optim pentru lungimea minimă de rutare definită în ecuația (4.15), se consideră următoarea problemă generalizată: Fiind date n puncte i ($i = 1, 2, \dots, n$) pe o linie, și un punct mobil x cu capacitatea c_{ix} între i și x , să se determine poziția $x = x_{opt}$ astfel încât

$$\sum_{i=1}^n c_{ix} |x - i| \quad (4.23)$$

să fie minim. Această problemă este un caz special al problemei de găsimă a unui centru absolut pentru un graf general. O condiție suficientă pentru determinarea punctului optim x_{opt} este dată de următoarea teoremă.

Teorema 4.4.1.1. Punctul x_{opt} este optim dacă este satisfăcută inegalitatea:

$$\sum_{i=1}^{x_{opt}-1} c_{ix} (x_{opt} - i) \leq \frac{N}{2} \leq \sum_{i=1}^{x_{opt}} c_{ix} (x_{opt} - i) \quad (4.24a)$$

unde:

$$N = \sum_{i=1}^n c_{ix} \quad (4.24b)$$

Revenind la ecuația (4.15), se consideră că partea dreaptă a acesteia este formată din doi termeni independenți. Lungimea totală de rutare pentru direcția x este:

$$D_x = \sum_{j=1}^m c_{Mj} |x - x_j| \quad (4.25)$$

Fie $x_{max} = \max_j \{x_j\}$, $x_{min} = \min_j \{x_j\}$, $j = 1, 2, \dots, m$. De notat că $x = x_{opt}$ trebuie să fie în domeniul:

$$x_{min} \leq x_{opt} \leq x_{max} \quad (4.26)$$

Se definește:

$$x(k) = x_j \quad (4.27)$$

$$c(k) = c_{Mj} \quad (4.28)$$

unde $x(k)$ și $c(k)$ sunt formele sortate ale lui x_j , respectiv c_{Mj} , iar $x(k)$ este în ordine ascendentă, $j = 1, 2, \dots, m$; $k = x_{min}, \dots, x_{max}$. Relația dintre j și k este determinată de procedura de sortare.

Din Teorema 4.4.1.1, o condiție suficientă pentru determinarea punctului optim x_{opt} este:

$$\sum_{k=x_{min}}^{x_{opt}-1} c(k)[x_{opt} - x(k)] \leq \frac{N}{2} \leq \sum_{k=x_{min}}^{x_{opt}} c(k)[x_{opt} - x(k)] \quad (4.29a)$$

unde:

$$N = \sum_{k=x_{min}}^{x_{max}} c(k) \quad (4.29b)$$

Similar, y_{opt} este dat de:

$$\sum_{k=x_{min}}^{y_{opt}-1} c(k)[y_{opt} - y(k)] \leq \frac{N}{2} \leq \sum_{k=x_{min}}^{y_{opt}} c(k)[y_{opt} - y(k)] \quad (4.30a)$$

unde:

$$N = \sum_{k=y_{min}}^{y_{max}} c(k) \quad (4.30b)$$

iar $y(k)$ este forma sortată a lui y_j în ordine ascendentă, cu:

$$y(k) = y_j \quad (4.31)$$

Relația dintre j și k este determinată de procedura de sortare pentru $j = 1, 2, \dots, m$ și $k = x_{min}, \dots, x_{max}$.

Dacă $x_{mopt} = x_{opt}$, $y_{mopt} = y_{opt}$, poziția (x_{mopt}, y_{mopt}) este numită *punct optim local* pentru modulul M_m în metrica Manhattan.

Până acum s-a indicat cum este selectat următorul modul candidat dacă există deja module plasate, și felul în care este plasat acest modul. Problema care se pune este selectarea și plasarea primului modul.

Dacă există module hard, ca module de I/E sau alte module specificate care sunt fixate în anumite poziții, acestea vor fi plasate primele. Dacă nu există module hard, există mai multe metode pentru selectarea primului modul. Se pot lua în considerare următoarele criterii:

1. Este selectat modulul cu gradul maxim, care este plasat în centru sau într-o altă poziție adecvată.
2. Este selectat modulul cu gradul minim, care este plasat într-un colț sau într-o altă poziție adecvată.

Următoarea problemă este plasarea modulului a cărei poziție optimă locală este ocupată de un alt modul fixat. În principiu, este posibilă încercarea tuturor pozițiilor.

țiilor disponibile și alegerea poziției celei mai bune. Acest proces este însă consumator de timp, și nu este avantajoasă utilizarea lui pentru plasarea inițială. Este necesară însă testarea pozițiilor din jurul poziției optime locale. Chiar dacă poziția optimă a modului candidat nu se află în jurul poziției optime locale, această metodă este mult mai rapidă și mai eficientă din cauza faptului că matricea de capacitate este rară.

Se prezintă în continuare un algoritm pentru plasarea constructivă inițială, elaborat de Wang și Chen [189].

1. Pe baza unui set dat S de seturi de semnale, se conectează fiecare set de semnale printr-un arbore.
2. Se determină matricea de capacitate.
3. Se testează dacă există module hard. În caz afirmativ, se plasează modulele hard în pozițiile cerute, apoi se continuă cu pasul 5. În caz contrar, se continuă cu pasul 4.
4. Se alege primul modul M_0 și se plasează în poziția selectată.
5. Se selectează un modul candidat.
6. Se determină punctul optim local W_{mopt} al modului candidat.
7. Se testează dacă W_{mopt} este ocupat. Dacă da, se testează setul care conține W_{mopt} și se plasează modulul în poziția optimă. În caz contrar, se plasează modulul în punctul optim local.
8. Se testează dacă M_m este ultimul modul. Dacă da, algoritmul se termină. În caz contrar, se continuă cu pasul 5.

Există diferite alte metode constructive. În secțiunea 4.4.5 va fi descrisă o metodă numerică numită plasare bazată pe forțe.

4.4.2 Plasarea pe baza tăieturii minime

În capitolul 3 a fost prezentată problema de partiționare a circuitelor. Un algoritm de partiționare încearcă gruparea modulelor puternic conectate. O asemenea grupare va reduce de asemenea lungimea interconexiunilor. Se va descrie în continuare un algoritm care utilizează în mod repetat o procedură de partiționare prin metoda tăieturii minime pentru generarea unei plasări.

În secțiunea 4.3 au fost prezentate trei funcții obiectiv, și anume $X(P)$, $Y(P)$ și $L(P)$. S-a arătat că prin minimizarea $X(P)$, tăietura orizontală maximă, și prin minimizarea $Y(P)$, tăietura verticală maximă, se va îmbunătăți rutabilitatea unei plasări pentru o rețea de porți.

Minimizarea funcției $X(P)$ este strâns legată de problema de bipartiționare. De aceea, se aplică un algoritm de partiționare pentru circuitul dat pentru a genera două blocuri A și B , se plasează modulele din blocul A la stânga unei linii imaginare de tăietură verticală c_1 , și se plasează modulele din blocul B la dreapta liniei de tăietură c_1 . Setul de tăietură obținut de algoritm este numărul de conexiuni orizontale tăiate de c_1 , și este notat cu $\Phi_P(c_1)$.

Presupunem că se repetă procesul pentru blocurile A și B , deci se consideră blocul A ca un circuit și se partiționează în două blocuri A_1 și A_2 , utilizând o linie de tăietură verticală c_2 . Similar, blocul B se partiționează în două blocuri B_1 și B_2 , utilizând o linie de tăietură verticală c_3 (Figura 4.5).

Acest proces poate fi repetat prin introducerea altor linii de tăietură. Presupunem că procedura de partiționare utilizată generează o partiție optimă. Pentru întregul circuit, $\Phi_P(c_1)$ este tăietura minimă posibilă. Similar, $\Phi_P(c_2)$ este tăietura minimă

posibilă pentru subcircuitul A . Nu se poate scrie că $\Phi_P(c_{i+1}) \leq \Phi_P(c_i)$, $1 \leq i \leq r-1$, deși în acest caz s-ar putea obține minimizarea funcției $X(P)$.

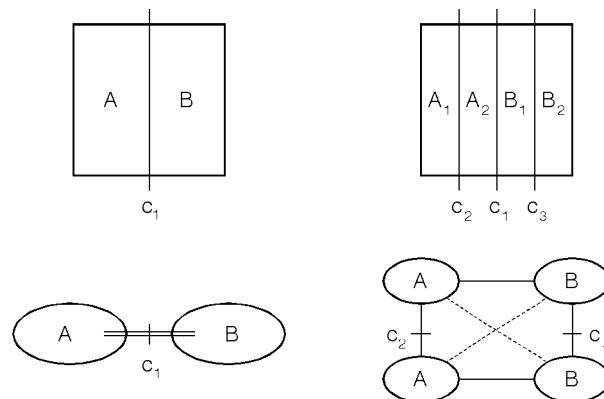


Figura 4.5. Utilizarea partiționării pentru a reduce $X(P)$.

Procedura descrisă mai sus nu minimizează $X(P)$, ci minimizează $\Phi_P(c_2)$ cu condiția ca $\Phi_P(c_1)$ să fie minim. Această funcție se poate scrie ca $\Phi_P(c_2) | \Phi_P(c_1)$. Procedura minimizează de asemenea $\Phi_P(c_3) | \Phi_P(c_1)$.

Minimizarea funcțiilor $X(P)$, $Y(P)$ sau $L(P)$ este foarte dificilă din punct de vedere computațional. Pentru simplificarea problemei, se utilizează o *funcție obiectivă secvențială*, notată cu $F(P)$, a cărei valoare apropiată de cea minimă este mai ușor de obținut.

$$F(P) = \min [\Phi_P(c_r)] | \min [\Phi_P(c_{r-1})] | \dots | \min [\Phi_P(c_1)] \quad (4.32)$$

unde c_1, c_2, \dots, c_r este o secvență ordonată de linii de tăietură verticale sau orizontale.

Algoritmul de plasare pe baza tăieturii minime presupune că este disponibilă o secvență ordonată de r linii de tăietură. Aceste r linii de tăietură împart suprafața de plasare în locații. Cerințele principale ale algoritmului sunt:

1. O procedură eficientă pentru partiționarea circuitului, și
2. O strategie de selecție a liniilor de tăietură.

Funcția $F(P)$ (ecuația 4.32) este minimizată prin partiționarea circuitului mai întâi în două, astfel încât numărul de conexiuni care traversează c_1 este minimizat. Dacă c_1 este o linie de tăietură verticală, atunci celulele din stânga liniei c_1 sunt fixate și nu se pot muta în dreapta. Celulele din dreapta liniei c_1 sunt de asemenea restricționate și nu se pot muta în stânga. Apoi, se utilizează următoarea linie de tăietură c_2 , și conexiunile care traversează c_2 sunt minimizează ținând cont de restricția impusă deja de c_1 . Procedura se continuă până când vor fi utilizate toate cele r linii de tăietură. Din cauza faptului că această procedură este de tip *greedy*, soluția obținută nu este în mod garantat cea optimă global.

Există trei metode mai utilizate pentru selecția liniilor de tăietură și a ordinii în care ele sunt prelucrate [146]. În cazul primei metode, numită *procedura de plasare quadratică*, suprafața de plasare este divizată în patru unități prin două linii de tăietură, una verticală și cealaltă orizontală, ambele trecând prin centru. Procedura de divizare anterioară este aplicată apoi recursiv fiecărui sfert din suprafață, până când întreaga suprafață este divizată în locații. Această secvență este ilustrată în Figura 4.6(a), fiind avantajoasă pentru circuite cu o densitate mare de rutare în centru.

În a doua metodă, numită *procedura de plasare prin biseccie*, suprafața este divizată în mod repetat în jumătăți egale prin linii de tăietură orizontale, obținându-se segmente orizontale. Această procedură de divizare este continuată până când fiecare segment orizontal este o linie, celulele fiind asignate acestor linii. Apoi, fiecare linie este biseccionată vertical în mod repetat, până când subregiunile rezultate conțin o locație. Această metodă este avantajoasă pentru plasarea celulelor standard și este ilustrată în Figura 4.6(b).

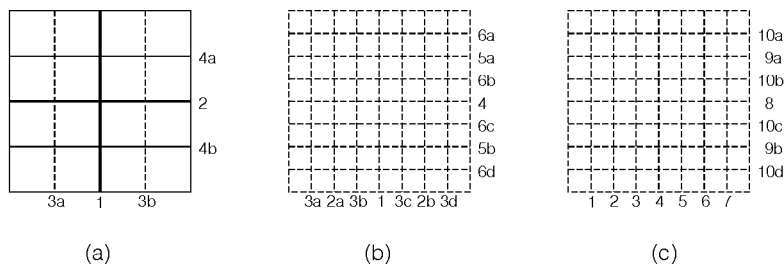


Figura 4.6. Ilustrarea secvențelor de aplicare a liniilor de tăietură.

În cazul ultimei metode, cele n celule ale circuitului sunt divizate prin linia de tăietură c_1 în două seturi de k și $(n - k)$ celule, astfel încât $\Phi_P(c_1)$ să fie minimizat. Primele k celule obținute sunt asignate liniei de sus (sau celei de jos). Procedura este aplicată apoi pentru cele $(n - k)$ celule, care sunt divizate în k și $(n - 2k)$ celule. Procesul este continuat până când toate celulele sunt asignate liniilor. Celulele sunt asignate apoi coloanelor utilizând biseccia verticală. Această secvență este ilustrată în Figura 4.6(c), fiind recomandată pentru celulele cu un număr mare de interconexiuni la periferie.

Structura unui algoritm recursiv de plasare pe baza tăieturii minime este prezentată în Figura 4.7.

```

Algorithm Tăietură_min ( $S, n, C$ );
    /*  $S$  este suprafața de plasare */
    /*  $n$  este numărul celulelor care trebuie plasate */
    /*  $n_0$  este numărul celulelor dintr-o locație */
    /*  $C$  este matricea de conectivitate */

begin
    if ( $n \leq n_0$ ) then plasare_celule ( $S, n, C$ )
    else
        ( $S_1, S_2$ )  $\leftarrow$  diviz_supr ( $S$ );
        ( $n_1, c_1$ ), ( $n_2, c_2$ )  $\leftarrow$  partit ( $n, C$ );
        Tăietură_min ( $S_1, n_1, c_1$ );
        Tăietură_min ( $S_2, n_2, c_2$ );
    endif
end.

```

Figura 4.7. Structura unui algoritm de plasare recursiv pe baza tăieturii minime.

Procedura de partiționare prezentată anterior nu ține cont de poziția pinilor de I/E (care este de obicei fixă) și de semnalele care intră într-un grup de celule. Aceste semnale afectează poziția în care trebuie plasate celulele în aceeași măsură în care afectează această poziție conexiunile din cadrul grupului. Includerea acestor semnale în procedura de plasare bazată pe partiționare se numește *propagarea terminalelor* [68].

4.4.3 Plasarea prin metoda călirii simulate

Metoda călirii simulate este una din cele mai utilizate pentru plasarea celulelor. Algoritmul general a fost prezentat în capitolul 3, secțiunea 3.4.4. În această secțiune se prezintă modul în care algoritmul poate fi adaptat pentru plasare. Se prezintă de asemenea algoritmul utilizat de pachetul de programe pentru plasare și rutare TimberWolf, dezvoltat de Carl Sechen și Sangiovanni-Vincentelli [151]. Acesta a fost primul program care a aplicat metoda de călire simulată pentru problema de plasare.

4.4.3.1 Aplicarea algoritmului de călire simulată pentru plasare

Algoritmul de călire simulată poate fi modificat pentru plasarea celulelor prin alegerea unei funcții adecvate de *perturbare* pentru a genera o nouă configurație de plasare, și prin definirea unei funcții corespunzătoare de *acceptare*. Pentru simplitate, se consideră că suprafața de plasare este modelată sub forma unei table de șah, și fiecare celulă a circuitului poate fi plasată într-un pătrat al tablei. O funcție simplă de vecinătate se obține prin interschimbarea perechilor, în care sunt alese două locații și se interschimă conținutul lor. Alte metode pentru generarea unor stări de vecinătate sunt mutarea unei celule selectate arbitrar într-o locație arbitrară, rotirea celulelor dacă aceasta este permisă de strategia de amplasare, sau orice altă mutare care poate modifica lungimea conexiunilor.

Fie $\Delta h = (Cost(NewS) - Cost(S))$ modificarea lungimii estimate a conexiunilor datorită unei interschimbări, unde $Cost(S)$ este vechea lungime a conexiunilor, iar $Cost(NewS)$ este lungimea după perturbare. Interschimbarea este acceptată dacă $\Delta h < 0$ (deci, $Cost(NewS) < Cost(S)$) sau dacă funcția de acceptare ($random < e^{-\Delta h / T}$) este adevărată, unde $random$ este un număr aleator între 0 și 1, generat în mod uniform, iar T este valoarea curentă a temperaturii.

4.4.3.2 Algoritmul TimberWolf

Pachetul de programe TimberWolf3.2 [151] este destinat configurațiilor de circuite cu celule standard. Pe baza datelor de intrare și a parametrilor furnizați de utilizator, programul de plasare construiește o topologie de circuit cu celule standard. Acești parametri, împreună cu lățimea totală a celulelor care trebuie plasate, permite programului să calculeze poziția inițială și lungimile rândurilor orizontale. Blocurile de macroui sunt plasate următoarele, urmate de plasarea celulelor de I/E. Blocurile de macroui și celulele de I/E își păstrează poziția inițială, fiind optimizată numai plasarea celulelor standard.

După plasarea inițială, algoritmul execută plasarea și rutarea în trei etape distincte. În prima etapă, celulele sunt plasate astfel încât să se minimizeze lungimea estimată a conexiunilor. În a doua etapă, sunt inserate celule de trecere după cum este necesar, lungimea conexiunilor este minimizată din nou, și se execută rutarea globală preliminară. În a treia etapă, sunt efectuate modificări locale ale plasării pentru a reduce numărul pistelor de rutare necesare. Se va prezenta în continuare prima etapă a algoritmului, care utilizează călirea simulată pentru plasare.

Funcția obiectiv care este minimizată de algoritmul TimberWolf3.2 în timpul plasării este costul estimat al interconexiunilor. Scopul primei etape este găsirea unei plasări a celulelor standard astfel încât costul total estimat al interconexiunilor să fie minimizat. Se utilizează o funcție de vecinătate numită *generate* pentru a produce noi stări prin efectuarea unei selecții aleatoare din trei *funcții de perturbare* posibile:

1. Mutarea unei singure celule într-o nouă locație, de exemplu într-un rând diferit.
2. Interschimbarea a două celule.
3. Oglindirea unei celule după axa x .

Perturbațiile sunt limitate la o regiune din cadrul unei ferestre de înălțime H_T și lățime W_T . De exemplu, dacă o celulă trebuie mutată, locația destinație este aleasă în cadrul unei ferestre de limitare centrată în jurul celei (Figura 4.8).

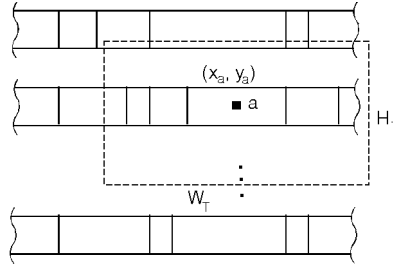


Figura 4.8. Fereastră de limitare centrată în jurul unei celule.

Două celule a și b , centrate în pozițiile (x_a, y_a) și (x_b, y_b) vor fi selectate pentru interschimbare numai dacă $|x_a - x_b| \leq W_T$ și $|y_a - y_b| \leq H_T$. Dimensiunile ferestrei sunt funcții descrescătoare cu temperatura T . Dacă temperatura curentă este T_1 și următoarea temperatură este T_2 , lățimea și înălțimea ferestrei sunt micșorate astfel:

$$W(T_2) = W(T_1) \frac{\log(T_2)}{\log(T_1)} \quad (4.33)$$

$$H(T_2) = H(T_1) \frac{\log(T_2)}{\log(T_1)} \quad (4.34)$$

Funcția de cost utilizată de algoritmul TimberWolf3.2 este suma a trei componente:

$$\gamma = \gamma_1 + \gamma_2 + \gamma_3 \quad (4.35)$$

Componenta γ_1 este o măsură a lungimii totale de interconectare estimate. Pentru fiecare conexiune i , dacă întinderea pe orizontală și pe verticală este X_i , respectiv Y_i , atunci lungimea estimată a conexiunii i este $(X_i + Y_i)$. Această valoare trebuie multiplicată cu ponderea w_i a conexiunii. Unei conexiuni i se pot atașa două ponderi, o componentă orizontală w_i^H și o componentă verticală w_i^V . Deci,

$$\gamma_1 = \sum_{i \in N} [w_i^H \cdot X_i + w_i^V \cdot Y_i] \quad (4.36)$$

unde însumarea se efectuează pentru toate conexiunile i din setul de conexiuni N .

Ponderea unei conexiuni este utilă pentru a indica măsura în care o conexiune este critică. Dacă o conexiune face parte dintr-o cale critică, de exemplu, trebuie ca aceasta să fie cât mai scurtă pentru a introduce o întârziere cât mai mică. Se pot mări ponderile conexiunilor critice pentru a se atinge acest scop. Ponderile independente pe orizontală și pe verticală permit utilizatorului flexibilitatea de a favoriza conexiunile într-o direcție față de cealaltă direcție.

Atunci când o celulă este mutată sau două celule sunt interschimbate, este posibil să existe o suprapunere între două sau mai multe celule. Fie O_{ij} suprafața de suprapunere între două celule i și j . Suprapunerile nu sunt de dorit și trebuie minimizate. A doua componentă a funcției de cost, γ_2 , este interpretată ca penalizare a suprapunerilor, și este definită astfel:

$$\gamma_2 = w_2 \sum_{i \neq j} [O_{ij}]^2 \quad (4.37)$$

În ecuația de sus, w_2 este ponderea penalizării. Prin ridicarea la pătrat a suprapunerii se asigură penalizări foarte mari pentru suprapuneri mai mari.

Datorită mutării celulelor și a interschimbării acestora, lungimea unui rând de celule poate deveni mai mare sau mai mică. A treia componentă a funcției de cost reprezintă o penalizare pentru cazul în care lungimea unui rând R depășește lungimea așteptată \overline{L}_R sau este mai scurtă decât aceasta:

$$\gamma_3 = w_3 \sum_{\text{rânduri}} |L_R - \overline{L}_R| \quad (4.38)$$

unde w_3 este ponderea diferenței de lungime. Distribuția inegală a lungimii rândurilor are ca efect irosirea spațiului. Există de asemenea o legătură între lungimea totală a conexiunilor și rutabilitate, pe de o parte, și distribuția inegală, pe de altă parte.

Funcția de răcire este reprezentată de

$$T_{i+1} = \alpha(T_i) \times T_i \quad (4.39)$$

unde $\alpha(T)$ este parametrul ratei de răcire, determinat experimental. Procesul de călire începe la o temperatură inițială foarte ridicată, de exemplu 4×10^6 . Inițial, temperatura este redusă rapid [$\alpha(T) \approx 0.8$]. În domeniul temperaturilor medii, temperatura este redusă lent [$\alpha(T) \approx 0.95$]. În domeniul temperaturilor joase, temperatura este redusă din nou rapid [$\alpha(T) \approx 0.8$]. Algoritmul se termină atunci când $T < 1$.

La fiecare temperatură, se încearcă efectuarea unui număr fix de mutări. Numărul optim al mutărilor depinde de dimensiunea circuitului. De exemplu, pentru un circuit cu 200 de celule, se recomandă 100 de mutări pe celulă, ceea ce presupune evaluarea a 2.34×10^6 configurații în aproximativ 125 de trepte de temperatură. Pentru un circuit cu 3000 de celule, se recomandă 700 de mutări pe celulă, ceea ce conduce la un număr total de 247.5×10^6 încercări [146].

4.4.4 Plasarea prin partiționare ierarhică

Deși metoda călirii simulate a fost aplicată cu succes pentru plasarea circuitelor, pe măsura creșterii dimensiunii circuitului timpul de calcul necesar devine inacceptabil în practică. Pentru reducerea timpului de calcul necesar algoritmului de călire simulată au fost utilizate diferite tehnici [88] [107] [119]. Hamada *et al.* [88] au propus utilizarea metodei de partiționare ierarhică pe baza tăieturii proporționale, elaborată de Cheng și Wei [50], urmată de aplicarea călirii simulate multinivel.

În prima etapă, circuitul este descompus într-un arbore de grupuri prin aplicarea recursivă a metodei de partiționare pe baza tăieturii proporționale. Se reduce astfel în mod semnificativ complexitatea problemei. Rezultatul partiționării este reprezentat printr-un arbore binar, a cărui rădăcină reprezintă întregul circuit. Fiecare nod intern reprezintă un subcircuit partiționat la un anumit nivel. În etapa a doua, fiecare grup este considerat ca o super-componentă, și se aplică metoda de călire simulată pentru acestea. În etapa a treia, se integrează soluțiile subproblemelor prin utilizarea metodei ferestrelor mobile. În fiecare pas, este definită o fereastră peste configurația curentă de plasare. Grupurile de sub această fereastră sunt repartiționate, și se aplică algoritmul de călire simulată noilor grupuri generate. Fereastra este apoi deplasată și ciclul continuă.

Algoritmul de plasare prin partiționare ierarhică este prezentat în Figura 4.9.

```

Algorithm Plasare_ierarhică;
/* Partiționare ierarhică și plasare */
begin
  Se calculează numărul de nivele necesare pentru partiționare nr_nivele;
  nivel = 0;
  for (nivel < nr_nivele - 1) do
    for (toate ferestrele de la acest nivel) do
      sel_fereastră (nivel, index_fereastră);
      Partiționare multinivel prin tăietura proporțională;
      TimberWolfSC;
      ajustare_rânduri ( $\rho$ );
      nivel ++;
    endfor
  endfor
/* Plasare finală */
for (toate ferestrele de la nivelul inferior) do
  Călire simulată la temperatură joasă prin TimberWolfSC;
endfor
end.

```

Figura 4.9. Algoritm de plasare prin partiționare ierarhică.

Prima fază a algoritmului este compusă din patru pași principali. În primul pas este selectată o fereastră din cadrul circuitului, pentru care configurația de plasare este rafinată în mai multe iterații. În al doilea pas, celulelor selectate li se aplică partiționarea ierarhică pe baza tăieturii proporționale, care generează grupuri în cadrul circuitului. În al treilea pas, grupurile sunt transformate în super-celule, pentru care se aplică programul de plasare *TimberWolfSC*. În ultimul pas, se efectuează ajustarea rândurilor pe care sunt plasate celulele, pentru a elimina distribuția neuniformă a grupurilor. În următoarea iterație a algoritmului, fereastra este deplasată în următoarea poziție, care se suprapune peste ferestrele procesate anterior de pe același nivel.

În ultima fază, cea de rafinare, grupurile sunt transformate în celule individuale, și sunt plasate utilizând călirea simulată la o temperatură joasă.

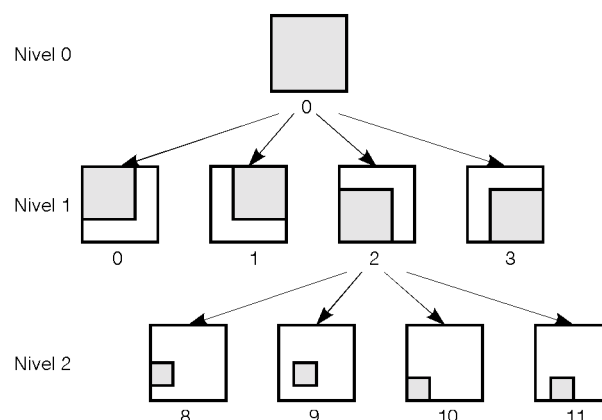


Figura 4.10. Plasare ierarhică prin metoda ferestrelor mobile.

Prima fereastră acoperă întregul circuit, și dimensiunea ferestrei se reduce pe măsură ce se ajunge la nivelele inferioare din ierarhie. Numărul de grupuri acoperite de o fereastră este menținut sub o anumită valoare constantă, numită capacitatea ferestrei. Poziția unui grup reprezintă poziția medie a celulelor din cadrul grupului.

În fiecare nivel, ferestrele sunt baleiate de la colțul din stânga sus al circuitului spre colțul din dreapta jos, de la un rând orizontal la altul. Figura 4.10 ilustrează utilizarea ferestrei în procesul de plasare. La nivelul 0, fereastra acoperă întregul circuit. La nivelul 1, fiecare din cele patru ferestre acoperă 43% din circuit, și ele se suprapun reciproc. La nivelul 2, o fereastră a nivelului 1 este acoperită de patru ferestre ale nivelului 2, deci există 16 ferestre la nivelul 2.

Acest concept al ferestrelor, combinat adesea cu o tehnică de partiționare succesivă, a fost utilizat de numeroase programe de plasare. De exemplu, programul PROUD, elaborat de Tsay și Kuh, utilizează biseția recursivă a circuitului în fiecare etapă de partiționare. Ceea ce este caracteristic pentru metoda utilizată de Hamada *et al.* [88] este faptul că circuitul este acoperit de o secvență de ferestre, în care două ferestre consecutive sunt suprapuse. Suprapunerea între ferestre permite trecerea celulelor peste limitele partițiilor, ceea ce poate compensa efectul partiției precedente în cazul în care aceasta a afectat calitatea plasării.

Ajustarea rândurilor cuprinde ca operație principală modificarea lungimii acestora prin mutarea celulelor dintr-un rând mai lung în rânduri adiacente mai scurte. Algoritmul procedurii de ajustare este prezentat în Figura 4.11.

```

Procedure ajustare_rânduri ( $\rho$ );
begin
  while (1) do
    Se selectează rândul cel mai lung rând_lung și un rând
    adiacent mai scurt rând_scurt;
    Se alege o celulă C din rând_lung și se calculează costul combinat:
       $\delta = \rho \cdot \Delta\text{lungime\_rând} + \Delta\text{lungime\_conexiune}$ ;
    if ( $\delta > 0$ ) then exit;
    else
      Se mută celula C din rând_lung în rând_scurt;
    endif
  endwhile
end

```

Figura 4.11. Procedura pentru ajustarea lungimii rândurilor.

Interschimbarea celulelor trebuie efectuată cu precauție, deoarece aceasta egalizează lungimea rândurilor în detrimentul lungimii totale a conexiunilor. Egalizarea lungimii rândurilor se termină atunci când minimul costului combinat

$$\delta = \rho \cdot \Delta\text{lungime_rând} + \Delta\text{lungime_conexiune}; \quad (4.40)$$

devine pozitiv. Primul termen al ecuației este negativ, deoarece trebuie minimizată diferența dintre lungimile rândurilor. Al doilea termen, care este lungimea conexiunilor sacrificată de egalizarea lungimii rândurilor, este pozitiv. Factorul ρ ponderează $\Delta\text{lungime_rând}$. Dacă ρ este setat la o valoare mare, lungimea rândurilor poate fi egalizată într-o mare măsură, cu prețul creșterii lungimii conexiunilor. În caz contrar, egalizarea lungimii rândurilor sacrifică doar o mică porțiune a lungimii conexiunilor.

În [88] s-a raportat pentru algoritmul de plasare prin partiționare ierarhică o reducere medie a lungimii totale a conexiunilor de 2.49% față de programul Timber-WolfSC, versiunea 5.6, și o reducere medie a timpului de calcul de 18.3% față de același program. S-a constatat de asemenea că timpul de calcul crește liniar cu dimensiunea circuitului [88].

4.4.5 Plasarea prin metode numerice

Problema de plasare poate fi transformată într-o problemă de optimizare numerică. În această secțiune se va prezenta o metodă numită plasare controlată de forțe, elaborată de Hanan și Kurtzberg [146]. Problema de plasare este redusă la soluționarea unui sistem de ecuații liniare pentru a determina locațiile de echilibru (coordonatele ideale x, y) ale celulelor.

Ideea de bază a acestei metode este că celulele interconectate exercită forțe unele asupra altora. Mărimea forței F exercitate de o celulă i asupra altei celule j este proporțională cu distanța care le separă. Aceasta este o analogie cu legea lui Hooke din mecanică, care se referă la forțele exercitate între două mase conectate printr-un arc. Dacă masele se află la o distanță d și constanta arcului este k , forța cu care masele se atrag este $k \times d$. Presupunem că o celulă a este conectată cu o altă celulă b printr-o conexiune cu ponderea w_{ab} . Fie d_{ab} distanța între a și b . Forța de atracție dintre celule este proporțională cu produsul $w_{ab} \times d_{ab}$. O celulă i conectată cu mai multe celule j aflate la distanțe d_{ij} prin conexiuni cu ponderi w_{ij} este atrasă cu o forță totală F_i dată de

$$F_i = \sum_j w_{ij} \cdot d_{ij} \quad (4.41)$$

Dacă celula i dintr-un asemenea sistem își poate modifica poziția, această deplasare se va realiza în direcția forței F_i până când forța rezultantă asupra celulei va fi zero. Locația în care se va deplasa celula este numită locație destinație de forță zero. În ecuația (4.41), F_i reprezintă lungimea totală ponderată a conexiunilor care pornesc de la celula i . Atunci când toate celulele se deplasează în locațiile lor de forță zero, suma pătratelor distanțelor este minimizată. Acesta este principiul care stă la baza metodei de plasare controlată de forțe.

Metoda constă în calcularea forțelor exercitate asupra oricărei celule date, și apoi deplasarea acesteia în direcția forței rezultante pentru plasarea celulei în locația sa de forță zero. Această locație (x_i^o, y_i^o) poate fi determinată prin egalarea cu zero a componentelor x și y ale forțelor exercitate asupra celulei, deci

$$\sum_j w_{ij} \cdot (x_j^o - x_i^o) = 0 \quad (4.42)$$

$$\sum_j w_{ij} \cdot (y_j^o - y_i^o) = 0 \quad (4.43)$$

Rezultă pentru x_i^o și y_i^o :

$$x_i^o = \frac{\sum_j w_{ij} \cdot x_j}{\sum_j w_{ij}} \quad (4.44)$$

$$y_i^o = \frac{\sum_j w_{ij} \cdot y_j}{\sum_j w_{ij}} \quad (4.45)$$

Trebuie evitată asignarea mai multor celule aceleiași locații, sau soluția banală care asignează toate celulele aceleiași locații (x_i^o, y_i^o).

Procedeu numeric pentru determinarea locațiilor ideale ale celulelor se poate extinde într-o procedură de plasare constructivă după cum urmează. Pornind de la o anumită plasare inițială, este selectată câte o celulă la un moment dat, și este calculată locația sa de forță zero. Procesul poate fi repetat în mod iterativ pentru a îmbunătăți soluția obținută. Deciziile care trebuie luate de un asemenea algoritm se referă la ordinea în care sunt selectate celulele și locația în care trebuie plasată celula în cazul în care locația de forță zero este ocupată.

Celula care va fi mutată poate fi selectată în mod aleator sau utilizând o tehnică euristică. Este logică selectarea celulei i pentru care F_i este maximă în configurația actuală. Dacă există restricții de timp, pot fi selectate mai întâi celulele care se află pe căile critice.

Dacă locația de formă zero a unei celule p este ocupată de o altă celulă q , sunt disponibile diferite opțiuni pentru plasarea celulei p , dintre care se prezintă următoarele:

- 1) Se mută p într-o locație liberă apropiată de q .
- 2) Se evaluează modificarea costului dacă p este interschimbabil cu q . Dacă rezultă o reducere a costului, se efectuează interschimbarea. Este necesară calcularea modificării costului deoarece este posibil ca celula q să se afle în locația sa de formă zero.
- 3) Celula p este plasată în locația ocupată, și este calculată o nouă locație de formă zero pentru celula q deplasată. Procedura este continuată până când sunt plasate toate celulele.
- 4) Celula p este plasată în locația calculată și celula q este mutată într-o locație adiacentă. Dacă locația adiacentă este ocupată de o altă celulă r , atunci celula r este mutată într-o locație adiacentă, și se continuă astfel până când se găsește o locație liberă.
- 5) Pentru evitarea acestei probleme se pot calcula locațiile de formă zero pentru toate celulele. Se caută perechi de celule de forma (p, q) astfel încât locația de formă zero a celulei p este poziția prezentă a celulei q și invers. Se interschimbă apoi celulele p și q , au fost găsite astfel locațiile de formă zero ale ambelor celule.

Există diferite versiuni ale metodei de plasare controlată de formă. Se prezintă în Figura 4.12 un algoritm cu îmbunătățire iterativă care utilizează mutarea de tip 3) pentru rezolvarea conflictului la o locație [153].

Algoritmii utilizează două indicatoare pentru fiecare locație a suprafeței de plasare. Un indicator OCUPAT arată dacă există o celulă asignată locației respective. Un indicator BLOCAT arată starea celulei asignate în mod curent locației. Dacă celula care ocupă locația a fost deplasată cel puțin o dată, indicatorul BLOCAT al locației va fi setat pentru a preveni deplasarea din nou a celulei respective.

În cadrul algoritmului este selectată câte o celulă, numită celulă sursă, fiind calculată conectivitatea și locația destinație de formă zero a acesteia. Selecția celulelor se bazează pe conectivitatea totală. Dacă locația destinație este ocupată, celula care ocupa în prealabil locația calculată va fi selectată pentru a fi mutată. Pentru a evita buclele infinite, oricare celulă care este mutată în locația sa destinație este blocată pentru valoarea curentă a contorului de iterații *contor_iter*. Buclele infinite pot apărea dacă două celule C_a și C_b concurează pentru aceeași locație destinație. Odată ce este selectată o celulă sursă și este calculată locația destinație de formă zero a acesteia, sunt posibile patru cazuri: locația calculată poate fi (1) *aceeași* cu locația inițială a celulei sursă, (2) o altă locație *liberă*, (3) o locație care este *ocupată* (dar *neblocată*), sau (4) o locație care este ocupată și *blocată*.

Bucula **while** interioară a algoritmului este executată cât timp *end_mutare* este *false*. Dacă locația calculată este aceeași cu cea curentă sau este o altă locație liberă, indicatorul *end_mutare* este setat la *true*, *contor_aband* este setat la zero și celula ocupă locația calculată. Este selectată apoi următoarea celulă sursă în ordinea conectivității și bucla continuă.

Dacă locația calculată este ocupată și nu este blocată, celula este mutată în locația calculată, iar celula care ocupa această locație este selectată ca următoarea celulă care va fi mutată. Indicatorul *end_mutare* este setat la *false*, iar *contor_aband* este setat la zero.


```

Algorithm Plasare_numerică;
begin
  Se calculează conectivitatea totală a fiecărei celule;
  Se sortează celulele în ordine descrescătoare a conectivității și se depun în lista L;
  while (contor_iter < limit_iter) do
    Sursă = următoarea celulă din L;
    Se declară poziția sursei ca liberă;
    while (end_mutare = false) do
      Se calculează locația destinație a sursei și se rotunjește la
      următorul întreg;
      case (locație_destinație) is
      LIBERĂ:
        Se mută celula sursă la destinație și se blochează;
        end_mutare = true;
        contor_aband = 0;
      ACEEAȘI CU LOCAȚIA CURENTĂ:
        end_mutare = true;
        contor_aband = 0;
      BLOCATĂ:
        Se mută celula selectată în locația liberă cea mai apropiată;
        end_mutare = true;
        contor_aband = contor_aband + 1;
        if (contor_aband > limit_aband) then
          Se deblochează toate locațiile celulelor;
          contor_iter = contor_iter + 1;
        endif;
      OCUPATĂ: /* și neblocată */
        Se selectează celula ca fiind următoarea care va fi mutată;
        Se mută celula sursă la destinație și se blochează;
        end_mutare = false;
        contor_aband = 0;
    endcase;
  endwhile;
endwhile;
end.

```

Figura 4.12. Exemplu de algoritm de plasare prin metode numerice.

În cazul în care locația destinație este ocupată și blocată, indicatorul *end_mutare* este setat la *true*, *contor_aband* este incrementat, iar celula este mutată în cea mai apropiată locație liberă. În acest caz, dacă *contor_aband* este mai mic decât *limit_aband*, este selectată următoarea celulă sursă, locațiile blocate rămân în aceeași poziție, și iterația continuă. Dacă însă *contor_aband* este mai mare decât valoarea prestabilită *limit_aband*, toate locațiile blocate sunt deblocate, este selectată o altă celulă sursă, contorul de iterații este incrementat și se începe o nouă iterație.

4.4.6 Plasarea liniară prin metode spectrale

Plasarea liniară este o problemă fundamentală pentru proiectarea circuitelor VLSI. O plasare liniară de calitate mai bună are ca efect reducerea lungimii conexiunilor, ceea ce reduce probabilitatea ca o conexiune să fie tăiată [139]. Au fost propuse diferite metode care utilizează vectori proprii, atât pentru problema de partiționare, cât și pentru plasarea liniară [39] [40] [86] [109].

În literatură s-au efectuat comparații între funcțiile obiectiv liniare și cuadractice. S-a constatat că prin utilizarea unei funcții liniare se obține o plasare de calitate mai bună din punct de vedere a lungimii conexiunilor. În [139], utilizarea unei funcții

liniare pentru plasare a permis de asemenea obținerea unor îmbunătățiri importante ale partiționării din punct de vedere a capacității tăieturii. Utilizarea unei funcții obiectiv cuadratică are însă ca efect obținerea unui număr mai redus de conexiuni foarte lungi față de cazul unei funcții obiectiv liniare. Cu alte cuvinte, deviația standard a lungimii conexiunilor este mai mică în cazul funcției cuadratică decât în cazul funcției liniare [109]. Aceasta înseamnă că funcția cuadratică are tendința să plaseze componentele într-un mod mai dispersat, rezultând mai puține componente suprapuse. Dezavantajul este însă că funcția cuadratică minimizează lungimea pătrată a conexiunilor în locul lungimii liniare, și de aceea nu corespunde direct cu scopul plasării liniare.

Li *et al.* [109] au propus utilizarea unei funcții obiectiv spectrale care este un compromis între funcția cuadratică și cea liniară, ceea ce permite folosirea avantajelor oferite de ambele funcții.

Metoda vectorilor proprii asigură nu numai o soluție euristică pentru problema de plasare liniară, dar furnizează și informații globale de conectivitate pentru gruparea componentelor. Aceasta permite reducerea dimensiunii problemei și poate îmbunătăți calitatea soluției obținute prin metodele euristice de partiționare și plasare. De asemenea, gruparea are avantaje pentru modelarea conexiunilor multi-pin. În urma grupării unui circuit, gradul hipermuchiiilor din graf este redus. În consecință, modelele grafurilor utilizate pentru aproximarea hipergrafului circuitului vor avea o acuratețe mai mare.

Problema de plasare liniară a unui circuit modelat printr-un hipergraf $H = (V_H, E_H)$ este de a se plasa componentele reprezentate prin setul de vârfuri V_H în locații aflate la distanțe egale (câte o componentă în fiecare locație), astfel încât lungimea totală a conexiunilor să fie minimă. În general, acest hipergraf este aproximat printr-un graf $G = (V, E)$, unde o hipermuchie este reprezentată printr-un set de muchii.

Fie $n = |V|$ numărul de vârfuri și $m = |E|$ numărul de muchii ale grafului G . Acest graf poate fi descris printr-o matrice de adiacență $A = [a_{ij}]$ cu dimensiunea $n \times n$, unde un element a_{ij} al matricii reprezintă ponderea conexiunii dintre vârfurile i și j . După transformarea unui hipergraf într-un graf, problema de plasare liniară se poate formula astfel:

$$\min \sum_{i>j} \sum_j a_{ij} |d_i - d_j| \quad (4.46)$$

unde d_i este coordonata locației pentru vârful i în cadrul plasării liniare.

Problema de plasare liniară este NP-completă. Printr-o metodă spectrală, plasarea spectrală continuă, în care restricția de plasare a vârfurilor în locații specifice este eliminată, se utilizează de obicei ca o euristică pentru soluționarea problemei de plasare liniară.

Fiind dat graful ponderat $G = (V, E)$, reprezentat prin matricea de adiacență $A = [a_{ij}]$, matricea Q , numită Laplacian al lui G , este definită astfel:

$$Q = \begin{cases} \sum_{j=1}^n a_{ij} & \text{dacă } i = j \\ -a_{ij} & \text{în caz contrar} \end{cases} \quad (4.47)$$

Cu metoda vectorilor proprii, problema de plasare liniară continuă este formulată ca o problemă de programare quadratică, astfel:

$$\min \sum_{i>j} \sum_j a_{ij} (x_i - x_j)^2 = X^T Q X \quad (4.48)$$

astfel încât

$$I^T X = 0, X^T X = 1$$

unde x_i este coordonata vârfului i în cadrul plasării liniare continue.

Prin calcularea celei mai mici valori proprii diferite de zero a Laplacianului Q și a vectorului propriu corespunzător X , se obține o soluție non-trivială a problemei de programare quadratică de mai sus, iar soluția euristică a problemei de plasare liniară se obține prin interpretarea vectorului propriu ca o ordonare liniară a vârfurilor V . Cu o asemenea metodă, funcția obiectiv care este minimizată este lungimea pătrată a conexiunilor.

În [139] a fost utilizată următoarea funcție obiectiv liniară pentru rezolvarea problemei de plasare liniară continuă:

$$\min \sum_{i>j} \sum_j a_{ij} |x_i - x_j| \quad (4.49)$$

astfel încât

$$\sum_i x_i = f$$

Prin utilizarea acestei funcții, s-a obținut o calitate ridicată a plasării și partiționării.

O asemenea funcție obiectiv liniară poate fi rescrisă ca o funcție quadratică prin modificarea a_{ij} cu distanța $|x_i - x_j|$ [109] [139]:

$$\sum_{i>j} \sum_j a_{ij} |x_i - x_j| = \sum_{i>j} \sum_j a'_{ij} (x_i - x_j)^2 \quad (4.50)$$

unde

$$a'_{ij} = \frac{a_{ij}}{|x_i - x_j|}$$

Astfel, prin îmbunătățire iterativă, metoda programării cuadractice poate fi utilizată pentru rezolvarea problemei de plasare liniară cu funcția de obiectiv liniară.

Pentru problema de plasare liniară, funcția obiectiv liniară permite obținerea unei calități mai ridicate a soluției decât funcția obiectiv quadratică; aceasta deoarece funcția liniară reprezintă o măsură mai exactă pentru problema de plasare liniară. Pe de altă parte, funcția obiectiv quadratică are avantajele amintite anterior.

Deoarece problema de plasare liniară este soluționată în mod euristic prin interpretarea vectorului propriu ca ordonare a vârfurilor, cu cât vârfurile sunt plasate mai dispersat, cu atât vor fi introduse mai puține erori numerice în cadrul plasării. De aceea, soluția continuă a plasării liniare trebuie să fie suficient de dispersată pentru a putea fi interpretată.

Pe baza acestor observații, Li *et al.* [109] au propus o funcție obiectiv de ordin α pentru problema de plasare liniară continuă, astfel:

$$\min \sum_{i>j} \sum_j a_{ij} |x_i - x_j|^\alpha = \sum_{i>j} \sum_j \frac{a_{ij}}{|x_i - x_j|^{2-\alpha}} (x_i - x_j)^2 \quad (4.51)$$

astfel încât

$$\sum_i x_i = f$$

unde $1.0 \leq \alpha \leq 2.0$. Atunci când $\alpha = 1.0$, funcția de ordin α devine funcția liniară, iar atunci când $\alpha = 2.0$, funcția de ordin α devine funcția quadratică. Cu această funcție obiectiv, se mărește dispersia soluțiilor, păstrându-se în același timp o măsură suficient de exactă pentru lungimea conexiunilor. Plasarea liniară continuă cu funcția obiectiv de ordin α poate fi soluționată în mod iterativ, și o asemenea metodă iterativă este

convergență, lungimea conexiunilor, definită în ecuația (4.51), fiind monoton descrescătoare [109].

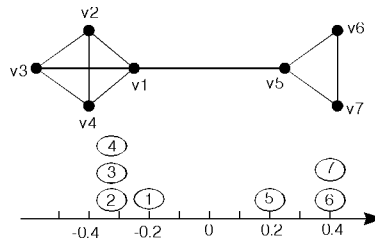


Figura 4.13. Exemplu de graf și plasarea liniară continuă prin metoda vectorilor proprii.

Soluția problemei de plasare liniară continuă nu reprezintă numai o euristică pentru soluționarea problemei de plasare liniară, dar și o euristică pentru problema de grupare a vârfurilor din graful circuitului. În Figura 4.13 se prezintă un graf simplu și plasarea liniară continuă obținută prin metoda vectorilor proprii. În acest exemplu, există 7 vârfuri $V = \{x_i \mid i = 1, \dots, 7\}$ și 10 conexiuni. Vectorul propriu corespunzător celei mai mici valori proprii diferite de zero $\lambda_1 = 0.3983$ pentru acest exemplu este $X = \{x_1 = -0.2142, x_2 = x_3 = x_4 = -0.3560, x_5 = 0.2965, x_6 = x_7 = 0.4928\}$.

Din această plasare liniară continuă, vârfurile V se pot divide în două grupuri. Primul grup cuprinde vârfurile ($v_i \mid i = 1, \dots, 4$), iar al doilea grup cuprinde vârfurile ($v_i \mid i = 5, 6, 7$). Aceasta deoarece vârfurile din fiecare grup sunt plasate foarte apropiat, în timp ce distanța dintre cele două grupuri este relativ mare. Conectivitatea din cadrul unui grup este deci mult mai ridicată decât conectivitatea între grupuri. Astfel, este naturală utilizarea rezultatelor plasării liniare continue ca și informații de conectivitate globală într-un proces de grupare.

Se prezintă în continuare un algoritm de plasare liniară care combină metoda spectrală cu gruparea circuitului [109]. Acest algoritm constă din trei faze: 1) prin utilizarea unei metode de grupare de jos în sus, se formează grupuri de noduri cu conectivitate ridicată; 2) se execută plasarea liniară a circuitului grupat; 3) se execută plasarea liniară a circuitului original. În ultima fază, se expandează mai întâi grupurile și se execută plasarea liniară pentru nodurile fiecărui grup. Apoi, se îmbunătățește calitatea plasării liniare printr-un algoritm care utilizează metoda fluxului maxim și a tăieturii minime. Algoritmul este prezentat în Figura 4.14.

Algoritmul de grupare a nodurilor circuitului folosește proprietatea de grupare a metodei vectorilor proprii ca și informație de conectivitate globală. Scopul principal al grupării este de a reduce dimensiunea circuitului fără a afecta calitatea partiționării și a plasării. Acest algoritm inițializează fiecare grup cu un singur nod al circuitului original. Se fuzionează apoi în mod iterativ perechi de grupuri care sunt puternic conectate, până când se reduce dimensiunea circuitului grupat până la o valoare predefinită.

Se utilizează următoarea metrică pentru a măsura conectivitatea unei perechi de grupuri i și j :

$$C = \frac{\sum_k c_{ij}^k}{W_i \times W_j} \times g(i, j) \quad (4.52)$$

unde c_{ij}^k reprezintă contribuția conexiunii cu numărul de ordine k la conectivitatea circuitului. Dacă grupurile i și j sunt conectate prin conexiunea k , această contribuție este

$$c_{ij}^k = \frac{1}{\dim(\text{conexiune}_k) - 1}$$

iar în caz contrar, $c_{ij}^k = 0$. W_i indică dimensiunea grupului i ; dacă dimensiunea grupului este mare, este descurajată participarea acesteia la operații suplimentare de grupare. În fine, $g(i, j)$ este ponderea rezultată din plasarea liniară continuă prin metoda vectorilor proprii, având rolul unei informații globale de conectivitate.

```

Algorithm Plasare_Liniară_Spectrală;
begin
1.   Se execută un algoritm de grupare a nodurilor circuitului;
     Fie  $\Theta = \{C_i \mid i = 1, \dots, K\}$  setul de grupuri;
2.   Se construiește noul circuit grupat pentru  $\Theta$ ;
     Se execută plasarea liniară cu funcția obiectiv de ordin  $\alpha$ ;
3.   Fie  $O = (C_1, \dots, C_K)$  ordonarea grupurilor;
      $i = 2$ ;
     while ( $i \leq K - 1$ ) do
         Se construiește pentru  $C_i$  o rețea rezistivă  $(s, v_{ji}, \dots, v_{jt}, t)$ ,
         unde  $v_{ji}$  ( $i = 1 \dots l$ ) este o componentă a lui  $C_i$ ;
         Se execută plasarea liniară a acestei rețele rezistive;
         Se înlocuiește  $C_i$  cu ordonarea  $(v_{i1}, \dots, v_{il})$  în  $O$ ;
          $i = i + 1$ ;
     endwhile;
end.

```

Figura 4.14. Algoritmul de plasare liniară prin metoda spectrală.

În ultima fază a algoritmului se utilizează conceptul rețelelor rezistive pentru rezolvarea problemei de plasare. Prin analogie cu un circuit electric, unde conductanța între nodul i și nodul j este egală cu $-q_{ij}$, elementul din matricea Q (Laplacian) definită în ecuația (4.47), s-a demonstrat că problema de plasare este echivalentă cu cea a algerii tensiunilor circuitului electric astfel încât puterea disipată să fie minimă. Tensiunile nodurilor circuitului electric reprezintă o analogie cu coordonatele nodurilor pentru problema de plasare. Spre deosebire de metoda vectorilor proprii, rețeaua rezistivă poate include specificațiile celulelor de I/E. Astfel, prin modelarea celulelor de I/E ca surse de tensiune fixă, se pot fixa anumite noduri la coordonate specifice, fiind necesară determinarea coordonatelor celorlalte noduri.

Metoda rețelelor rezistive poate fi utilizată pentru determinarea coordonatelor nodurilor dintr-un grup, rezultând decompimarea circuitului grupat. Fiind dată ordonarea grupurilor $O = (C_1, C_2, \dots, C_K)$, atunci când se realizează decompimarea grupului C_i , se încearcă determinarea coordonatelor nodurilor din C_i în timp ce ordonarea grupurilor nu este schimbată. Prin metoda rețelelor rezistive, se comprimă toate grupurile de la stânga grupului C_i din O într-o celulă de I/E s cu tensiunea 0, se înlocuiește grupul C_i cu componentele sale, și se comprimă toate grupurile de la dreapta grupului C_i din O într-o celulă de I/E t cu tensiunea 1. Apoi se determină coordonatele nodurilor din C_i pe baza tensiunilor acestei rețele rezistive comprimate. Pentru aceasta se utilizează următoarea teoremă:

Teorema 4.4.6. [109] În rețeaua rezistivă comprimată de sus, tensiunea unui nod mobil este între 0 și 1.

Pe baza acestei teoreme, nodurile din C_i vor fi plasate între două noduri comprimate s și t . Astfel, procesul de decompimare nu va modifica ordonarea grupurilor definită de O .

După decompimarea grupului $C_i \mid i = 2, \dots, K-1$ în ordonarea O , se obține o soluție a problemei de plasare liniară pentru circuitul original. Această metodă de gru-

pare și decompunere permite obținerea unei plasări de calitate mai ridicată din punct de vedere a lungimii conexiunilor decât metoda vectorilor proprii cu o funcție obiectiv liniară.

4.4.7 Alte metode iterative

4.4.7.1 Plasarea prin asignare cuadratică

Problema de asignare cuadratică poate fi definită astfel: Fiind dată o matrice a costurilor $\mathbf{C} = [c_{ij}]$, o matrice a distanțelor $\mathbf{D} = [d_{ij}]$, și o permutare p a primilor n întregi, să se determine minimumul expresiei

$$Q = \sum_{i,j} c_{ij} d_{p(i)p(j)} \quad (4.53)$$

pentru toate permutările p .

Una din interpretările posibile ale acestei formulări este următoarea: Fiind date n persoane (module) și o matrice $\mathbf{C} = [c_{ij}]$, unde c_{ij} este o măsură a afinității (capacității) între persoanele i și j (modulele i și j), n birouri posibile (locații) pentru aceste persoane (module), și o matrice $\mathbf{D} = [d_{st}]$, unde d_{st} este distanța măsurată cu o anumită metrică între birourile (locațiile) s și t . Dacă i este asignat biroului (locației) $p(i)$, iar j este asignat biroului (locației) $p(j)$, costul acestei asignări este $c_{ij} d_{p(i)p(j)}$. Astfel, costul total pentru fiecare i și j este Q , definit în ecuația (4.52). Se presupune că afinitatea reprezintă mărimea comunicării directe. Obiectivul este de a minimiza distanța totală parcursă (lungimea totală de rutare) de toate persoanele.

Deoarece există $n!$ permutări pentru n întregi, există $n!$ posibilități de asignare a modulelor la cele n locații.

Comparând problema de asignare cuadratică cu problema generală de plasare, se observă că problema de asignare cuadratică este un caz special al problemei de plasare, prin presupunerea că toate seturile de semnale conțin numai două module. În cazul problemei de asignare cuadratică, se consideră numai perechi de puncte, în timp ce în cazul problemei de plasare se consideră un set de puncte sau seturi de semnale. Totuși, tehnicile pentru rezolvarea problemelor de asignare cuadratică pot fi utilizate pentru problemele de plasare dacă se utilizează următoarea transformare: se setează c_{ij} egal cu suma ponderilor seturilor de semnale comune cu modulele i și j pentru $i \neq j$, cu $c_{ij} = 0$, și se alege matricea $\mathbf{D} = [d_{ij}]$ ca fiind una măsurată într-o anumită metrică a problemei originale de plasare.

O soluție optimă pentru problema asociată de asignare cuadratică nu este în mod necesar optimă pentru plasarea originală datorită transformării matematice utilizate. De aceea, transformarea utilizată afectează în mare măsură rezultatele.

4.4.7.2 Plasarea pe baza optimizării rețelelor rezistive

Ideea de utilizare a analogiei cu rețelele rezistive pentru soluționarea problemelor de plasare a fost introdusă mai întâi de Charney și Plato, iar apoi de Cheng și Kuh, care au prezentat o metodă sistematică pentru soluționarea problemei de plasare prin această analogie [189].

Funcția obiectiv este dată de:

$$\Phi(x, y) = \frac{1}{2} \sum_{i,j=1}^n c_{ij} d_{ij}^2 = \frac{1}{2} \sum_{i,j=1}^n c_{ij} [(x_i - x_j)^2 + (y_i - y_j)^2] \quad (4.54)$$

unde d_{ij} este distanța euclidiană între modulele i și j , iar x și y reprezintă lungimile conexiunilor măsurate în direcția x , respectiv y . Într-un mod mai concis, ecuația de sus se poate scrie ca:

$$\Phi = \mathbf{X}^T \mathbf{Q} \mathbf{X} + \mathbf{Y}^T \mathbf{Q} \mathbf{Y} \quad (4.55)$$

unde $\mathbf{Q} = \mathbf{G} - \mathbf{C}$ este o matrice simetrică cu dimensiunea $n \times n$, $\mathbf{C} = [c_{ij}]$ este matricea de capacitate, iar \mathbf{G} este o matrice diagonală cu $g_{ij} = \sum_{j=1}^n c_{ij}$ [hal]. \mathbf{X} și \mathbf{Y} sunt vectorii distanțelor măsurate în direcția x , respectiv y . Funcția obiectiv corespunde puterii P disipate într-o rețea rezistivă:

$$P = \mathbf{V}_n^T \mathbf{Y}_n \mathbf{V}_n \quad (4.56)$$

unde \mathbf{V}_n este vectorul reprezentând tensiunile nodurilor, iar \mathbf{Y}_n este matricea simetrică a admitanțelor.

Astfel, problema de plasare a modulelor poate fi formulată ca o problemă de optimizare liniară a unei rețele rezistive, dacă se efectuează următoarea transformare: modulele fixe corespund nodurilor cu tensiuni constante, iar modulele care pot fi mutate corespund nodurilor ale căror tensiuni trebuie determinate.

4.4.7.3 Plasarea prin algoritmul Steinberg

Algoritmul Steinberg folosește un set independent de module în timpul fiecărei iterații, utilizând o procedură optimă pentru reconfigurarea acelui set [189]. Ideea de bază este că atunci când este ales un set independent, modulele din acel set sunt eliminate de pe suprafața de plasare, și se calculează costul plasării fiecăruia din aceste module în fiecare locație disponibilă, obținându-se o soluție optimă. Deoarece aceste costuri sunt independente de locațiile celorlalte module din setul independent, poziționarea acestor module devine o problemă de asignare liniară. După se execută iterațiile pentru toate seturile independente care sunt generate, ciclul se termină și începe următorul ciclu.

În general, procesul de construire a unui set independent începe prin alegerea aleatoare a unui modul, și continuă prin adăugarea la acest set doar a acelor module care nu sunt conectate cu membrii setului. Acest proces continuă până când s-a format un set cu dimensiunea dorită sau nu mai pot fi adăugate alte module la set. Dimensiunea setului independent care va fi utilizat depinde în mare măsură de puterea de calcul disponibilă pentru problema de plasare. Este de preferat utilizarea seturilor independente maxime ori de câte ori este posibil pentru a se asigura flexibilitatea maximă în cadrul problemei de plasare.

O îmbunătățire a algoritmului Steinberg a fost sugerată de Rutman, în scopul obținerii unor rezultate mai bune și pentru a obține convergența mai rapidă a algoritmului.

4.4.7.4 Plasarea prin metoda spațiului grafurilor

Unii autori au încercat îmbunătățirea plasării utilizând metode bazate pe teoria grafurilor. Una din aceste metode este numită metoda spațiului grafurilor, în care relațiile de conexiune între modulele din cadrul unui circuit electric sunt reprezentate printr-un hipergraf. O muchie a unui hipergraf poate conecta orice număr de noduri, în timp ce o muchie a unui graf obișnuit nu poate conecta mai mult de două module. Prin utilizarea acestei metode, nodurile (reprezentând module) sunt mapate în spațiul grafului astfel încât distanța între nodurile din spațiu reflectă ponderea (numărul de conexiuni) unei muchii între nodurile hipergrafului original. Se minimizează apoi lungimea totală de rutare prin aplicarea tehnicii primei asignări propuse de Hung și Rom [189].

O altă metodă care utilizează un *digraf relaxat* a fost propusă de Ciesielski și Kinnen pentru soluționarea problemei de plasare bidimensională cu scopul de a se minimiza spațiul utilizat de module cu dimensiuni arbitrare.

4.4.7.5 Plasarea bazată pe operații cu linii și coloane

Se consideră o placă bidimensională pe care se plasează module în pozițiile definite prin coordonatele x și y . Fiecare modul va aparține unei linii și coloane. Un algoritm bazat pe operații cu linii și coloane constă din doi pași. Întâi, se determină o plasare optimă locală pentru module într-o linie (coloană), și apoi se obține plasarea îmbunătățită între două linii (coloane). Procesul este iterat și se termină atunci când se obține un rezultat satisfăcător sau s-a epuizat un anumit timp de calcul.

Problema de determinare a plasării optime într-o linie (coloană) este definită astfel: Fiind date n module, n locații într-o linie, și o matrice de capacitate $\mathbf{C} = [c_{ij}]$, unde c_{ij} este ponderea între modulele i și j , trebuie determinată lungimea minimă de rutare prin asignarea unui singur modul fiecărei locații. Numărul de soluții existente este $n!$, astfel că pentru un număr mare de module nu se pot examina toate soluțiile posibile. Pentru problema asociată de asignare, metoda elaborată de Munkres poate găsi soluția într-un timp rezonabil, în cazul unui număr de module de ordinul 200. Un algoritm euristic a fost elaborat de Wang și Chen pe baza noțiunii de arbore maxim de acoperire într-un graf general.

4.5 Metode neconvenționale de plasare

În această secțiune sunt prezentate unele metode neconvenționale utilizate pentru rezolvarea problemei de plasare. În secțiunea 4.5.1 este prezentată implementarea paralelă a unui algoritm de plasare prin metoda călirii simulate. În secțiunea 4.5.2 se prezintă utilizarea rețelelor neuronale artificiale pentru plasare.

4.5.1 Plasarea prin algoritmi paraleli

În această secțiune se prezintă implementarea paralelă a unui algoritm de plasare care se bazează pe metoda de călire simulată. Implementarea paralelă a metodei de călire simulată nu este simplă, din cauza naturii secvențiale a acestei metode. Călirea simulată poate fi descrisă ca o secvență de lanțuri Markov omogene [142]: fiecare pas de calcul al unui lanț începe doar atunci când pasul precedent s-a terminat. Aceasta este condiția ca întregul proces să conducă la o configurație fezabilă unică. Pentru a reduce timpul de calcul, se pot utiliza două tipuri de paralelism:

- Un paralelism pentru evaluarea fiecărei mutări: calculul unui anumit pas al lanțului Markov depinde numai de configurația sistemului înaintea acestui pas și este executat fără interacțiune cu ceilalți pași. Astfel, evaluările diferitelor mutări pot fi executate în paralel, ca și calculele variației funcției de cost și a criteriului de acceptare. Acest tip de paralelism este dependent de problemă.
- Un paralelism global la nivelul lanțului Markov, care poate fi combinat cu primul tip de paralelism, dacă este necesar.

În literatură au fost raportate diferite implementări paralele ale metodei de călire simulată [104] [141]. Diferențele constau în principal în următoarele aspecte:

- condițiile de convergență ale algoritmului paralel;
- dependența paralelismului de problema care trebuie rezolvată.

Pentru implementarea paralelă a problemei de plasare au fost sugerate diferite soluții. Metoda utilizată de Casotto *et al.* [37] constă în partiționarea setului de celule care trebuie plasate într-un număr de subseturi egal cu numărul procesoarelor disponibile; fiecare subset fiind asignat unui anumit procesor. Procesoarele funcționează asincron cât timp mutările apar într-un set dat de celule. Sunt permise interschimbările celulelor între diferite seturi, dar aceasta presupune suspendarea unuia din cele două procesoare implicate. Mutările din cadrul unui set de celule sunt translații, rotații sau interschimbări între celule. Deoarece procesoarele funcționează asincron, niciunul din procesoare nu cunoaște exact care este configurația curentă. Această metodă a fost implementată pe un calculator Sequent Balance 8000 cu maxim 8 procesoare. Casotto a experimentat de asemenea paralelizarea masivă a algoritmului de călire simulată pe calculatorul Connection Machine.

Mallela și Grover [119] au utilizat de asemenea grupuri de celule pentru a reduce numărul celulelor care trebuie plasate în cadrul fiecărei subprobleme. Plasarea celulelor în fiecare grup implică un spațiu de căutare redus, deci un timp de calcul redus, fiind posibilă evaluarea în paralel a fiecărui grup.

O altă abordare a problemei de plasare a celulelor a fost sugerată de Damera, Kirkpatrick și Norton. În acest caz, fiecare procesor evaluează o perturbație a lanțului Markov, cu condiția că două procesoare nu pot muta aceleași celule simultan; astfel, nu există conflict între procesoare și configurația finală este întotdeauna validă [142]. Atunci când o perturbație este acceptată, configurația celulelor este actualizată, indiferent de mutările care se calculează. La temperaturi joase, când rata de acceptare este redusă, nu există o diferență importantă față de algoritmul secvențial. La temperaturi mai înalte, comportarea metodei paralele diferă în mod semnificativ de cea a metodei secvențiale. Măsurătorile au fost efectuate pe un mediu de emulare care permite simularea unui sistem multiprocesor cu memorie partajată cu până la 64 de procesoare.

O metodă diferită a fost utilizată de Kravitz și Rutenbar [104], care au introdus noțiunea de subset de mutări serializabile. Un set de mutări este serializabil dacă mutările nu interacționează unele cu altele. Dacă mutările unui asemenea subset sunt evaluate în paralel, rezultatul este același ca și în cazul evaluării secvențiale. Determinarea acestor subseturi este însă din ce în ce mai dificilă și consumatoare de timp pe măsură ce numărul de procesoare crește. Soluția care a fost sugerată la această problemă este considerarea "celui mai simplu subset serializabil", unde este acceptată o singură mutare acceptabilă, în timp ce toate mutările respinse sunt contorizate. La temperaturi înalte, algoritmul paralel este foarte diferit de cel secvențial, deoarece rata de acceptare nu este aceeași.

Roussel-Ragot și Dreyfus [142] au sugerat o implementare paralelă a algoritmului de călire simulată care, pe de o parte, este independentă de problemă, iar pe de altă parte, are aceleași proprietăți de convergență ca și algoritmul serial. Aceștia au utilizat două moduri de paralelizare, în funcție de valoarea temperaturii, și au elaborat modele statistice care pot estima creșterea de viteză pentru orice problemă, în funcție de rata de acceptare și numărul de procesoare.

Rata de acceptare $\chi(T)$ este raportul între numărul de mutări acceptate și numărul de mutări încercate, pentru o anumită temperatură. Una din caracteristicile algoritmului de călire simulată este reducerea ratei $\chi(T)$ pe măsură ce temperatura scade. Cele mai multe metode paralele de călire simulată utilizează acest fapt. În regimul de temperaturi joase, când rata de acceptare χ este redusă, se pot utiliza K procesoare astfel încât $K < 1/\chi$. Deci, va fi acceptată cel mult o mutare în timp ce sunt evaluate K mutări, astfel încât este de așteptat ca timpul de calcul în modul paralel să fie mai redus decât timpul de calcul în modul secvențial cu un factor de ordinul K .

Pe lângă timpul de calcul, convergența algoritmului este de importanță deosebită în scopul obținerii unei soluții valide (sau a uneia optime). Unele metode de paralelizare pot afecta convergența algoritmului, în unele cazuri această convergență

fiind chiar imposibilă. Au fost efectuate diferite studii teoretice ale algoritmului secvențial, dar pentru algoritmul paralel de călire simulată teoria este mult mai puțin dezvoltată. De aceea, este avantajoasă utilizarea unei metode paralele care respectă condițiile de convergență ale algoritmului secvențial.

În cazul metodei propuse în [142], pentru ca algoritmul paralel să aibă o convergență similară cu cel secvențial, se urmărește generarea unui lanț Markov ale cărui stări au aceeași distribuție a probabilității ca și lanțul secvențial, dacă toți ceilalți parametri sunt aceeași. În funcție de rata de acceptare $\chi(T)$, se consideră două regimuri:

- Un regim al temperaturilor joase: Dacă $\chi(T) < 1/K$, va fi acceptată cel mult o mutare din K . Procesoarele încearcă mutări pe cont propriu, în mod asincron, în paralel, până când unul din cele K procesoare acceptă o mutare. Dacă este găsită o mutare acceptată, procesoarele sunt sincronizate, memoriile lor sunt actualizate cu noua configurație, și începe următorul pas de evaluare.
- Un regim al temperaturilor înalte ($\chi(T) > 1/K$): În acest regim, fiecare procesor poate evalua o singură mutare, și așteaptă până când toate celelalte procesoare își termină evaluarea. Apoi, se alege în mod aleator una din mutările acceptate, memoriile procesoarelor sunt actualizate cu noua configurație, și începe următorul pas de evaluare.

În plus pe lângă cele K procesoare "slave" menționate, metoda necesită un procesor "master", care monitorizează procesul de călire, alege mutarea acceptată în regimul temperaturilor înalte, actualizează memoria fiecărui procesor, și ține evidența statisticilor.

Pentru elaborarea unui model statistic al celor două regimuri de temperatură, se notează cu:

L_a	numărul maxim de mutări acceptate la o temperatură dată;
L_t	numărul maxim de încercări la o temperatură dată;
τ_0	timpul mediu de calcul necesar pentru evaluarea unei mutări în modul secvențial;
τ_r	timpul mediu necesar comunicației cu cele K procesoare "slave" și sincronizării acestora.

Semnificația pentru τ_0 este clară numai dacă se utilizează un singur tip de mutare elementară; în caz contrar, τ_0 reprezintă valoarea medie a timpului de calcul pentru diferitele tipuri de mutări care pot apare în timpul călirii simulate. În ambele cazuri, τ_0 poate fi estimat ca durata unui pas de temperatură împărțită cu numărul de mutări încercate. Dacă perturbațiile variază cu temperatura (de exemplu, dacă interschimbările blocurilor apar în principal la temperaturi înalte și translațiile apar mai ales la temperaturi joase), valoarea τ_0 va fi variabilă cu temperatura.

Temperatura este micșorată fie când numărul de mutări acceptate la temperatura curentă ajunge la L_a , fie când numărul de mutări încercate la temperatura curentă ajunge la L_t , indiferent care din aceste limite este atinsă prima.

A. Regimul temperaturilor înalte. În acest regim, fiecare procesor evaluează o mutare, și toate procesoarele sunt sincronizate la sfârșitul fiecărei evaluări. Timpul mediu necesar pentru execuția unei evaluări de către cele K procesoare va fi $\tau_0 + \tau_r$. De notat că τ_r ține cont de faptul că durata mutărilor poate fi diferită, astfel încât timpul necesar terminării unei evaluări paralele este egal cu durata celei mai lungi mutări încercate.

Deoarece lungimea lanoului Markov depinde de numărul de mutări acceptate și/sau de numărul de mutări încercate, trebuie evaluate mai întâi aceste cantități. Se presupune că, după o evaluare paralelă a K mutări, r mutări din K sunt respinse. Deci, $K - r$ mutări vor fi acceptabile, dar numai una din ele va fi acceptată în lanoul Markov. Raportul dintre numărul de mutări acceptate și numărul de mutări încercate, în modul secvențial, este $(K - r) / K$. În modul paralel, o singură mutare va fi acceptată.

Se pune problema construirii unui lanou cu aceeași rată de acceptare ca și în modul secvențial. Pentru aceasta, se poate proceda în felul următor: Se numerotează procesoarele într-o ordine arbitrară, de la 1 la K . Se notează cu n numărul primului procesor din listă care acceptă o mutare, și se construiește lanoul Markov cu primele n mutări din listă. În modul secvențial, un lanou Markov cu $n - 1$ mutări respinse, urmat de o mutare acceptată, ar fi avut aceeași probabilitate. Se poate arăta că valoarea medie pentru n este egală cu $n^* = (K + 1) / (K - r + 1)$ [142].

Dacă toate cele K procesoare resping mutările lor încercate, numărul primului procesor nu poate fi determinat utilizând această metodă, dar deoarece nu există mutare acceptată, numărul mutărilor încercate care trebuie luate în considerare este K .

Deci, atunci când se execută în paralel K evaluări, se poate scrie că:

- Dacă cel puțin o mutare a fost acceptată, se alege în mod aleator una din aceste mutări, și se consideră că au fost încercate un număr de mutări egal cu $n^* = (K + 1) / (K - r + 1)$.
- Dacă nici o mutare nu a fost acceptată, se consideră că au fost încercate K mutări.

Aceasta permite o estimare corectă a numărului de mutări încercate efectiv, deoarece se poate arăta că, în medie, raportul dintre numărul de mutări acceptate și numărul efectiv de mutări încercate în modul paralel este egal cu χ , rata de acceptare în modul secvențial [142]. Astfel se obține aceeași convergență în modul paralel ca și în modul secvențial, deoarece algoritmul paralel are aceeași matrice de probabilitate a tranzițiilor ca și cel secvențial.

Se evaluează în continuare numărul total efectiv de mutări încercate N_t^* , și numărul total de mutări acceptate N_a^* , după ce au fost executate N evaluări paralele pentru K mutări.

Numărul mediu efectiv de încercări este dat de [142]

$$N_t^* = N \cdot \left[(1 - \chi)^K \cdot K + \sum_{i=1}^K \binom{K}{i} \cdot \chi^i (1 - \chi)^{K-i} \frac{K+1}{i+1} \right] \quad (4.57)$$

sau, în mod echivalent,

$$N_t^* = N \cdot \frac{1 - (1 - \chi)^K}{\chi} \quad (4.58)$$

Limita L_t a mutărilor încercate este atinsă după un număr N_t de evaluări paralele a K mutări, care este dat de

$$N_t = L_t \cdot \frac{\chi}{1 - (1 - \chi)^K} \quad (4.59)$$

Numărul N_a^* de mutări acceptate este egal cu numărul de evaluări paralele a K mutări care conduc la cel puțin o mutare acceptabilă, deci

$$N_a^* = N[1 - (1 - \chi)^K] \quad (4.60)$$

Astfel, limita L_a a mutărilor acceptate este atinsă după un număr N_a de evaluări paralele a K mutări, care este dat de

$$N_a = \frac{L_a}{1 - (1 - \chi)^K} \quad (4.61)$$

Deci, în regimul temperaturilor înalte, numărul de evaluări paralele la o anumită temperatură este

$$N_p = \min(N_t, N_a) \quad (4.62)$$

Timpul de calcul corespunzător este

$$t_p = N_p (\tau_0 + \tau_r) \quad (4.63)$$

În modul serial, timpul de calcul este

$$t_s = \tau_0 \cdot L_a / \chi, \text{ dacă limita } L_a \text{ este atinsă prima}$$

$$t_s = \tau_0 \cdot L_t, \text{ dacă limita } L_t \text{ este atinsă prima}$$

Astfel, oricare din limitele L_a sau L_t este atinsă prima

$$\frac{t_p}{t_s} = \frac{\chi}{1 - (1 - \chi)^K} \left(1 + \frac{\tau_r}{\tau_0} \right) \quad (4.64)$$

De notat că $\lim (t_p/t_s) = 1 + \tau_r/\tau_0$ atunci când $\chi \rightarrow 1$, iar $\lim (t_p/t_s) = 1/K \cdot (1 + \tau_r/\tau_0)$ atunci când $\chi \rightarrow 0$.

La temperaturi înalte, eficiența este scăzută deoarece este respins un număr mic de mutări, rata de acceptare fiind ridicată. La temperaturi joase, timpul de calcul este divizat de numărul de procesoare, dacă timpul de comunicație τ_r este redus comparativ cu τ_0 .

Valoarea medie pentru τ_0 este cunoscută de la algoritmul secvențial. Determinarea valorii τ_r nu este simplă și depinde de problemă. Dacă τ_0 este constant, τ_r poate fi aproximat cu timpul de comunicație multiplicat cu K : dacă cele K calcule în paralel se termină în același timp, vor fi necesare K comunicații succesive pentru ca procesorul "master" să cunoască toate rezultatele și să relanseze operațiile executate de procesoarele "slave".

B. Regimul temperaturilor joase. În acest regim, fiecare procesor evaluează mutările în mod independent până când unul din cele K procesoare acceptă o mutare. La sfârșitul fiecărei evaluări individuale, procesoarele transmit în mod asincron procesorului "master" rezultatul evaluării. Dacă nici o mutare nu a fost acceptată de la comunicația precedentă, se încearcă o altă mutare. Dacă a fost acceptată o mutare, memoriile procesoarelor "slave" sunt actualizate și procesoarele sunt sincronizate. În acest mod, toate mutările respinse sunt contorizate ca pași către realizarea stării de echilibru.

Pentru a modela comportarea la temperatura joasă, este necesară evaluarea numărului de mutări necesare pentru ca o mutare să fie acceptată. Pentru aceasta, se poate estima timpul necesar pentru mutările individuale pe baza numărului de mutări în cazul algoritmului secvențial. În medie, o configurație este acceptată atunci când sunt evaluate $1/\chi$ mutări. În modul serial, procesul a avansat deci cu $1/\chi$ pași către starea de echilibru. Deoarece se dorește obținerea unei configurații fezabile a sistemului, dacă o altă mutare este acceptată de una din celelalte $K - 1$ procesoare, aceasta nu este luată în considerare. Perturbațiile sunt alese în mod aleator, și dacă, de exemplu, unul din blocuri este interschimbat în două mutări diferite, configurația rezultată nu este validă, deoarece blocul ar trebui plasat în două locații diferite. Atunci când procesoarele sunt sincronizate, au fost evaluate $1/\chi + K - 1$ mutări și, în medie,

au fost acceptate $(1/\chi + K - 1) \cdot \chi$ mutări. Deoarece sunt eliminate toate mutările acceptate, cu excepția uneia, se contorizează în lanțul Markov doar un număr de pași egal cu:

$$(1/\chi + K - 1) - (1/\chi + K - 1) \cdot \chi + 1 = 1/\chi + (K - 1) \cdot (1 - \chi)$$

Dacă τ_m este timpul necesar pentru a obține o mutare acceptată în paralel, comportarea în regimul temperaturilor joase poate fi modelat astfel:

i. Dacă limita L_a este atinsă prima:

$$t_p = L_a \cdot \tau_m \quad (4.65)$$

$$t_s = L_a / \chi \cdot \tau_0 \quad (4.66)$$

deci

$$\frac{t_p}{t_s} = \chi \frac{\tau_m}{\tau_0} \quad (4.67)$$

ii. Dacă limita L_t este atinsă prima:

$$t_p = \frac{L_t}{1/\chi + (K - 1)(1 - \chi)} \tau_m \quad (4.68)$$

$$t_s = L_t \cdot \tau_0 \quad (4.69)$$

deci

$$\frac{t_p}{t_s} = \frac{1}{1/\chi + (K - 1)(1 - \chi)} \frac{\tau_m}{\tau_0} \quad (4.70)$$

De observat că pentru $K = 1$, $\tau_m = \tau_0 / \chi$, astfel încât $t_p = t_s$.

Valoarea medie pentru τ_0 este cunoscută, dar determinarea valorii τ_m nu este simplă și depinde de problemă. Dacă τ_0 este constant și dacă $1/\chi$ este mult mai mare decât K , valoarea lui τ_m poate fi aproximată prin $(\tau_0 + \tau_c) / K\chi$, unde τ_c este timpul necesar unui procesor "slave" pentru comunicarea rezultatului.

Roussel-Ragot și Dreyfus au utilizat algoritmul paralel de călire simulată pentru o problemă simplă de plasare a unei rețele bidimensionale de b^2 circuite pe o suprafață pătrată. În configurația inițială a sistemului, fiecare circuit este conectat cu cei mai apropiați vecini prin conexiuni cu două terminale. Mutarea elementară este interschimbarea a două circuite, alese în mod arbitrar. Funcția de cost este lungimea totală a conexiunilor. Parametrii utilizați sunt următorii:

- Configurația inițială este aleasă în mod aleator.
- Temperatura inițială este aleasă astfel încât rata de acceptare este de peste 0.9.
- Temperatura este modificată atunci când au fost evaluate $5 \cdot b^2 \cdot (b^2 - 1)$ mutări sau atunci când au fost acceptate $b^2 \cdot (b^2 - 1) / 2$ mutări.
- Parametrul de răcire α este egal cu 0.9.
- Procesul de călire simulată se oprește atunci când temperatura atinge valoarea 0.2 sau când nu mai este acceptată nici o mutare la o temperatură dată.

Algoritmul a fost implementat pe o rețea de transputere. Un transputer a fost utilizat ca procesor "master", calculele fiind efectuate pe transputerile conectate cu acest procesor.

4.5.2 Plasarea prin rețele neuronale artificiale

Interesul recent manifestat pentru rețelele neuronale are la bază recunoașterea faptului că creierul uman execută calcule într-un mod diferit de cel al calculatoarelor digitale. Calculatoarele pot fi extrem de rapide și precise în executarea secvențelor de instrucțiuni care au fost formulate în mod explicit. Un sistem uman de procesare a informațiilor este format din neuroni cu o viteză de comutare de aproximativ un milion de ori mai redusă decât cea a porților logice. Cu toate acestea, un sistem uman este mult mai eficient decât calculatoarele în rezolvarea unor probleme complexe din punct de vedere computațional, ca de exemplu înțelegerea vorbirii. De asemenea, omul poate procesa informații vizuale într-un mod mai eficient decât cele mai rapide calculatoare.

Rețelele neuronale artificiale au fost utilizate pentru rezolvarea unor probleme dificile, inclusiv pentru unele probleme de optimizare. În secțiunea 4.5.2.1 se prezintă unele concepte fundamentale și modele ale rețelelor neuronale artificiale. În secțiunea 4.5.2.2 se prezintă rețelele neuronale Hopfield care se pot utiliza pentru rezolvarea diferitelor probleme de optimizare, inclusiv a celei de plasare. În secțiunea 4.5.2.3 se prezintă modul în care se pot utiliza rețelele neuronale artificiale pentru rezolvarea problemei de plasare.

4.5.2.1 Concepte de bază ale rețelelor neuronale artificiale

O rețea neuronală artificială poate fi definită ca o rețea sintetică care emulează rețelele neuronale biologice ale organismelor vii. O altă definiție este că rețelele neuronale artificiale reprezintă o clasă de algoritmi matematici, deoarece o rețea poate fi considerată ca o notație grafică pentru o clasă largă de algoritmi [190].

O rețea neuronală artificială este un ansamblu al unui mare număr de *neuroni artificiali*. Prima definiție formală a unui model de neuron artificial bazată pe o considerare foarte simplificată a unui neuron biologic a fost formulată de McCulloch și Pitts. Acest model este ilustrat în Figura 4.15.

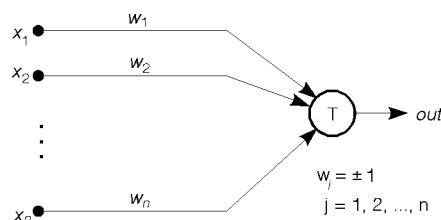


Figura 4.15. Modelul McCulloch-Pitts al neuronului.

Intrările x_i , pentru $i = 1, 2, \dots, n$, sunt 0 sau 1, depinzând de absența sau prezența impulsului de intrare la momentul de timp k . Semnalul de ieșire al neuronului este out . Regula de activare a acestui model este definită astfel:

$$out^{k+1} = \begin{cases} 1 & \text{dacă } \sum_{i=1}^n w_i x_i^k \geq T \\ 0 & \text{dacă } \sum_{i=1}^n w_i x_i^k < T \end{cases} \quad (4.71)$$

unde indicii superiori $k = 0, 1, 2, \dots$ indică momentul discret de timp, iar w_i este ponderea conexiunii intrării i cu membrana neuronului. T este valoarea de prag a neuro-

nului, care trebuie depășită de suma ponderată a intrărilor pentru ca neuronul să fie activat. Deși acest model al neuronului este foarte simplist, are totuși un important potențial computațional. Modelul poate executa operațiile logice de bază NOT, OR și AND, dacă valorile ponderilor și ale pragurilor sunt selectate în mod corespunzător.

Modelul McCulloch-Pitts al unui neuron este caracterizat prin formalism și o definiție matematică precisă. Modelul utilizează însă unele simplificări importante: permite numai stări binare, operează cu presupunerea unui timp discret, și presupune sincronismul funcționării tuturor neuronilor dintr-o rețea de dimensiuni mari. Ponderile intrărilor și valorile de prag ale neuronilor sunt fixe. Totuși, acest model se poate utiliza pentru introducerea unor concepte de bază.

Un neuron artificial constă dintr-un *element de procesare (nod)*, care recepționează un număr de intrări analogice având conexiuni sinaptice, și generează o singură ieșire. Simbolul general al unui neuron este indicat în Figura 4.16. Semnalele de intrare x_i ale neuronului sunt considerate unidirecționale, ca și semnalul de ieșire *out*. Fiecărei intrări x_i se asociază o pondere w_i .

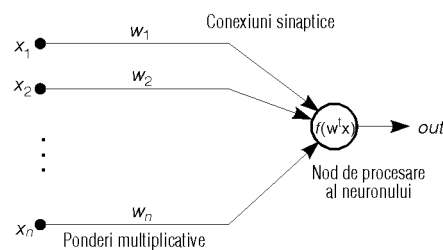


Figura 4.16. Simbolul general al neuronului constând din nodul de prelucrare și conexiunile sinaptice.

Semnalul de ieșire al neuronului este dat de relația următoare:

$$out = f(\mathbf{w}^T \mathbf{x}), \text{ sau} \quad (4.72)$$

$$out = f\left(\sum_{i=1}^n w_i x_i\right) \quad (4.73)$$

unde \mathbf{w} este vectorul ponderilor definit ca

$$\mathbf{w} = [w_1 \ w_2 \ \dots \ w_n]^T$$

iar \mathbf{x} este vectorul de intrare:

$$\mathbf{x} = [x_1 \ x_2 \ \dots \ x_n]^T$$

Ambii vectori sunt vectori coloană. Funcția $f(\mathbf{w}^T \mathbf{x})$ este numită *funcție de activare* a neuronului. Domeniul acestei funcții este setul valorilor de activare, *net*, ale modelului neuronului, de aceea această funcție este utilizată adesea ca $f(\text{net})$. Variabila *net* este definită ca un produs scalar al vectorului ponderilor și al vectorului de intrare:

$$net = \mathbf{w}^T \mathbf{x} \quad (4.74)$$

Argumentul funcției de activare, variabila *net*, este analogul potențialului membranei neuronului biologic. De observat că valoarea de prag T nu este utilizată în mod explicit în ecuațiile (4.72), (4.73) și (4.74). Deoarece în unele modele valoarea de prag are un rol important, în aceste cazuri este necesară extragerea explicită a acestei valori ca un parametru separat al modelului neuronului. Până acum s-a presupus că neuronul modelat are $n-1$ conexiuni sinaptice de la variabilele de intrare x_1, x_2, \dots, x_{n-1} . S-a presupus de asemenea că $x_n = -1$ și $w_n = T$.

Diferitele clase de rețele neuronale artificiale utilizează diferite funcții $f(net)$. De asemenea, chiar în cazul aceleiași clase, neuronii se pot comporta diferit în timpul diferitelor faze ale funcționării rețelei. De aceea, modelul general al neuronului este înlocuit de obicei cu un model specific, pentru o anumită funcție $f(net)$.

Se observă din ecuațiile (4.72) și (4.73) că neuronul, ca un nod de procesare, execută produsul scalar sau operația de însumare a intrărilor sale ponderate pentru a obține variabila net . Ulterior, execută operația neliniară $f(net)$ prin funcția sa de activare. Funcții de activare tipice sunt [190]

$$f(net) = \frac{2}{1 + e^{-\lambda net}} - 1 \quad (4.75)$$

și

$$f(net) = \text{sgn}(net) = \begin{cases} +1, & net > 0 \\ -1, & net < 0 \end{cases} \quad (4.76)$$

unde $\lambda > 0$ din ecuația (4.75) determină panta funcției continue $f(net)$ în apropiere de $net = 0$. Pentru $\lambda \rightarrow \infty$, limita funcției continue devine funcția $\text{sgn}(net)$ definită în ecuația (4.76). Funcțiile de activare (4.75) și (4.76) reprezintă funcția *bipolară continuă*, respectiv *bipolară binară*. Cuvântul "bipolar" indică faptul că se generează atât răspunsuri pozitive, cât și răspunsuri negative ale neuronilor pentru aceste definiții ale funcției de activare.

Prin deplasarea și scalarea funcțiilor de activare definite de ecuațiile (4.75) și (4.76) se poate obține funcția de activare *unipolară continuă*, respectiv *unipolară binară*, definite prin:

$$f(net) = \frac{1}{1 + e^{-\lambda net}} \quad (4.77)$$

și

$$f(net) = \begin{cases} 1, & net > 0 \\ 0, & net < 0 \end{cases} \quad (4.78)$$

Funcțiile de activare (4.75) și (4.77) se numesc adesea *caracteristici sigmoidale*, spre deosebire de funcțiile de activare (4.76) și (4.78), care descriu modelul discret al neuronului.

Dacă funcția de activare a neuronului are forma bipolară binară din ecuația (4.76), simbolul din Figura 4.16 se poate înlocui prin diagrama din Figura 4.17(a), care este o schemă bloc funcțională a unui neuron discret, indicând însumarea efectuată de nodul de sumare și limitarea efectuată de unitatea logică de prag.

În cazul funcției de activare continue cu forma din ecuația (4.75), modelul utilizat este ilustrat în Figura 4.17(b). Neuronul este reprezentat ca un amplificator sumator care amplifică semnalul de intrare $\mathbf{w}^T \mathbf{x}$. Modelele din Figura 4.17(a) și (b) se numesc *perceptron discret (binar)*, respectiv *perceptron continuu*. Perceptronul discret, introdus de Rosenblatt, a fost prima mașină de învățare, fiind precursorul multor modele neuronale actuale.

Pe baza definiției neuronului artificial, o rețea neuronală artificială poate fi definită ca o interconexiune de neuroni astfel încât ieșirile neuronilor sunt conectate, prin ponderi, cu ceilalți neuroni; sunt permise conexiuni fără întârzieri sau cu întârzieri. Există două modele principale ale rețelelor neuronale artificiale: rețele directe (*feedforward*) și rețele cu reacție inversă (*feedback*).

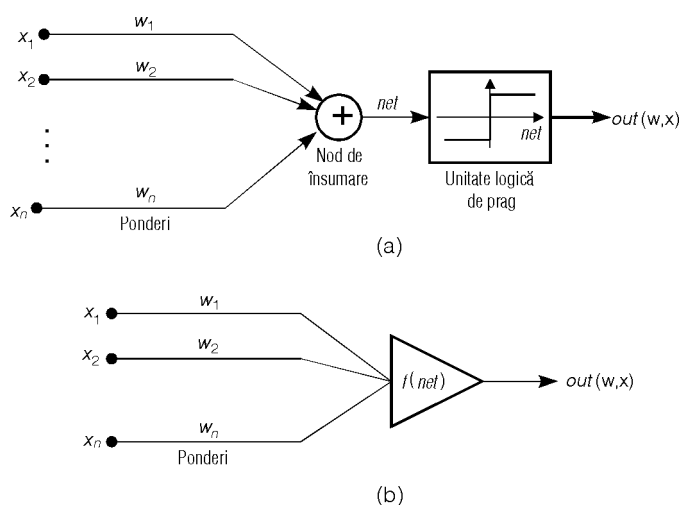


Figura 4.17. Modele ale neuronilor cu conexiuni sinaptice: (a) perceptron discret (binar); (b) perceptron continuu.

O rețea neuronală directă formată din m neuroni, fiecare având n intrări, este prezentată în Figura 4.18. Vectorii de intrare și de ieșire ai rețelei sunt:

$$\mathbf{x} = [x_1 \ x_2 \ \dots \ x_n]^T \tag{4.79}$$

$$\mathbf{out} = [out_1 \ out_2 \ \dots \ out_m]^T \tag{4.80}$$

Conexiunea cu ponderea w_{ij} conectează neuronul i cu intrarea j . Valoarea de activare pentru neuronul i se poate scrie ca:

$$net_i = \sum_{j=1}^n w_{ij}x_j, \text{ pentru } i = 1, 2, \dots, m \tag{4.81}$$

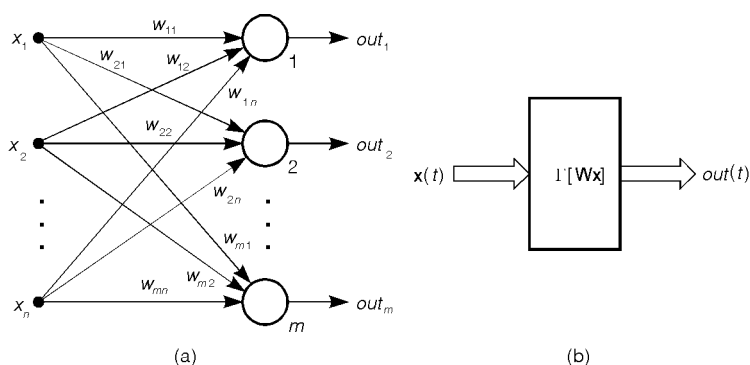


Figura 4.18. Rețea feedforward cu un singur strat: (a) schema de interconectare; (b) schema-bloc.

Următoarea transformare neliniară este efectuată de fiecare din cei m neuroni din rețea asupra vectorului de intrare \mathbf{x} :

$$out_i = f(\mathbf{w}_i^T \mathbf{x}), \text{ pentru } i = 1, 2, \dots, m \tag{4.82}$$

unde vectorul ponderilor \mathbf{w}_j conține ponderile conexiunilor care conduc spre nodul de ieșire i , și este definit astfel:

$$\mathbf{w}_i = [w_{i1} \ w_{i2} \ \dots \ w_{in}]^T \tag{4.83}$$

Dacă se introduce operatorul matricial neliniar Γ , maparea spațiului intrărilor \mathbf{x} în spațiul ieșirilor \mathbf{out} implementată de rețea se poate exprima astfel:

$$\mathbf{out} = \Gamma[\mathbf{W} \mathbf{x}] \tag{4.84}$$

unde \mathbf{W} este *matricea ponderilor*, numită de asemenea *matricea de conexiuni*:

$$\mathbf{W} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1n} \\ w_{21} & w_{22} & \dots & w_{2n} \\ \vdots & \vdots & \dots & \vdots \\ w_{m1} & w_{m2} & \dots & w_{mn} \end{bmatrix} \tag{4.85}$$

iar

$$\Gamma[\cdot] = \begin{bmatrix} f(\cdot) & 0 & \dots & 0 \\ 0 & f(\cdot) & \dots & 0 \\ \vdots & \vdots & \dots & \vdots \\ 0 & 0 & \dots & f(\cdot) \end{bmatrix} \tag{4.86}$$

Funcțiile neliniare de activare $f(\cdot)$ de pe diagonala operatorului matricial Γ operează asupra valorilor de activare *net* ale fiecărui neuron. Fiecare valoare de activare este un produs scalar a unei intrări cu vectorul ponderilor.

Rețele neuronale directe sunt caracterizate prin lipsa reacției inverse. Aceste rețele pot fi conectate în cascadă pentru a crea o rețea multistrat. Într-o asemenea rețea, ieșirea unui strat reprezintă intrarea pentru stratul următor.

O rețea neuronală cu reacție inversă poate fi obținută din rețeaua directă din Figura 4.18 prin conectarea ieșirilor neuronilor la intrările lor. Rezultatul este ilustrat în Figura 4.19.

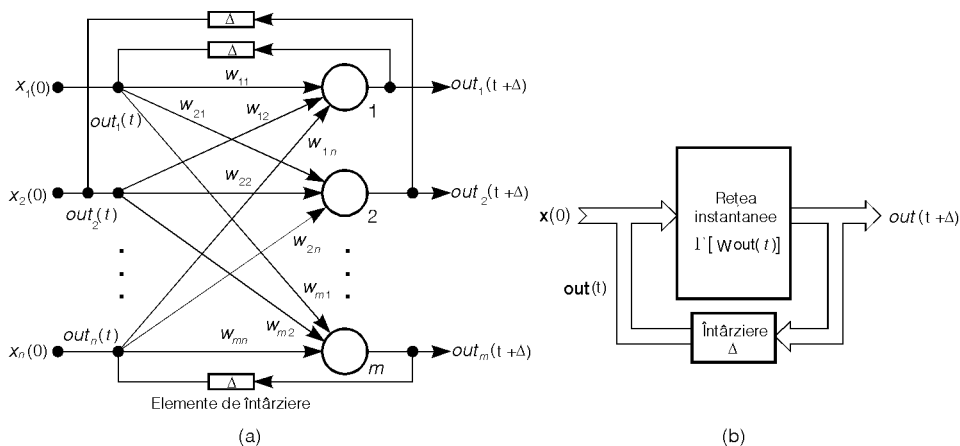


Figura 4.19. Rețea cu reacție inversă cu un singur strat: (a) schema de interconectare; (b) schema-bloc.

Rolul închiderii buclei de reacție este de a se permite controlul ieșirii out_i prin ieșirile out_j , pentru $j = 1, 2, \dots, m$. Un asemenea control este util dacă ieșirea prezentă, de exemplu $out(t)$, controlează ieșirea de la următorul moment de timp, $out(t + \Delta)$.

Timpul Δ dintre momentele t și $t + \Delta$ este introdus prin elementele de întârziere din bucla de reacție. Maparea ieșirii $\text{out}(t)$ în $\text{out}(t+\Delta)$ se poate scrie ca:

$$\text{out}(t+\Delta) = \Gamma[\mathbf{W} \cdot \text{out}(t)] \quad (4.87)$$

Intrarea $\mathbf{x}(t)$ este necesară doar pentru a inițializa rețeaua astfel încât $\text{out}(0) = \mathbf{x}(0)$. Intrarea este apoi eliminată și sistemul rămâne autonom pentru $t > 0$. Aici se consideră deci un caz special al acestei configurații cu reacție, astfel încât $\mathbf{x}(t) = \mathbf{x}(0)$, și nu este furnizată nici o intrare pentru $t > 0$.

Există două categorii principale de rețele cu reacție inversă. Dacă se consideră timpul ca o variabilă discretă, performanțele rețelei fiind observate la momente de timp discrete $\Delta, 2\Delta, 3\Delta, \dots$, sistemul este numit *cu timp discret*. Pasul de timp este considerat unitar, momentele de timp fiind indexate prin întregi pozitivi. Simbolul Δ are atunci semnificația întârzierii unitare. Pentru o rețea neuronală cu timp discret, ecuația (4.87) se transformă în

$$\text{out}^{k+1} = \Gamma[\mathbf{W} \cdot \text{out}^k], \text{ pentru } k = 1, 2, \dots \quad (4.88)$$

unde k este momentul de timp.

Rețeaua din Figura 4.19 este numită *recurentă* deoarece răspunsul acesteia la momentul de timp $k+1$ depinde de întreaga istorie a rețelei începând de la $k = 0$. Rețelele recurente operează de obicei cu o reprezentare discretă a datelor. Un sistem cu intrări la momente discrete de timp și cu o reprezentare discretă a datelor este numit un automat. Rețelele neuronale recurente din această clasă pot fi considerate deci ca automate [190].

Ecuația (4.88) descrie starea out^k a rețelei la momentele $k = 1, 2, \dots$, și generează o secvență de tranziții ale stărilor. Rețeaua începe tranzițiile stărilor după ce este inițializată la momentul 0 cu \mathbf{x}^0 , trecând printr-o succesiune de stări până când, eventual, ajunge într-o stare de echilibru. Această stare de echilibru este numită adesea *atractor*. Un atractor poate consta dintr-o singură stare sau un număr limitat de stări. Secvența de stări a unei rețele recurente este în general non-deterministică. În plus, pot exista mai multe stări de echilibru care pot fi atinse în mod potențial de rețea după o secvență de asemenea tranziții non-deterministice.

4.5.2.2 Rețele neuronale Hopfield

Rețelele neuronale Hopfield se pot utiliza pentru rezolvarea diferitelor probleme de optimizare, printre care și problema de plasare, motiv pentru care sunt descrise pe scurt în această secțiune. Pentru asemenea probleme, rețelele neuronale de tip gradient sunt cele mai adecvate. La aceste rețele, funcția de energie descrește în timp, timpul fiind presupus o variabilă continuă.

O rețea neuronală Hopfield este caracterizată ca o rețea puternic interconectată de procesoare simple analogice. O asemenea rețea de tip gradient converge spre unul din minimele stabile din spațiul stărilor [114] [190]. Evoluția sistemului este în direcția generală a gradientului negativ a unei funcții de energie. În mod tipic, funcția de energie a rețelei este echivalată cu o anumită funcție obiectiv (penalizare) care trebuie minimizată. Căutarea unei energii minime executată de o rețea de tip gradient corespunde căutării unei soluții a unei probleme de optimizare.

Pentru rezolvarea unor probleme de tipul celei a comis-voiajorului, Hopfield și Tank au definit o funcție de energie quadratică:

$$E(\mathbf{V}) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N T_{ij} V_i V_j - \sum_{i=1}^N V_i I_i \quad (4.89)$$

unde $\mathbf{V} = (V_1, \dots, V_N)$ reprezintă starea rețelei neuronale. V_i este ieșirea neuronului i . T_{ij} este intensitatea conexiunii sinaptice între neuronul V_i și V_j , cu $T_{ii} = 0$. I_i reprezintă

curentul de intrare. Rețeaua neuronală necesită conexiuni simetrice pentru a asigura convergența numai în stările de echilibru.

Pentru rezolvarea problemelor de programare liniară, Tank și Hopfield au definit o funcție de energie liniară:

$$E(\mathbf{V}) = \mathbf{A}^T \mathbf{V} + \sum_{j=1}^m F(\mathbf{W}_j^T \mathbf{V} - b_j) \quad (4.90)$$

unde $\mathbf{A}^T \mathbf{V}$ este funcția de cost supusă condiției restrictive $\mathbf{W}_j^T \mathbf{V} \geq b_j$. F este o funcție neliniară a cărei derivată f este definită ca

$$f(x) = \begin{cases} 0, & x \geq 0 \\ x, & x < 0 \end{cases}$$

Fiecare minim stabil este un atractor în spațiul stărilor. Setul stărilor inițiale care inițiază evoluția terminată într-un atractor este numit *rezervor de atracție*. Atractorul și rezervorul fiecărui atractor sunt determinate de funcția de energie și de parametrii interni ai rețelei neuronale. O rețea neuronală Hopfield reprezintă o procedură de căutare locală deterministică. După selectarea stării inițiale, rețeaua va converge la starea minimă în rezervorul căreia se află starea inițială [114].

Ca un exemplu simplu de rețea neuronală Hopfield, presupunem că valoarea analogică de intrare x trebuie convertită în reprezentarea sa binară printr-o rețea cu doi neuroni. Se presupune că funcția de activare este continuă unipolară, fiind definită între 0 și 1. Eroarea de conversie A/D E_c poate fi considerată ca o funcție de energie, care poate fi exprimată prin următoarea relație:

$$E_c = \frac{1}{2} \left(x - \sum_{i=0}^1 V_i 2^i \right)^2 \quad (4.91)$$

unde $V_0 + 2V_1$ este valoarea zecimală a vectorului binar $[V_1 \ V_0]^T$, și corespunde valorii analogice x . Motivul acestei alegeri a funcției de energie este că minimizarea energiei va minimiza simultan și eroarea de conversie. Scopul este determinarea matricii \mathbf{W} și a vectorului curenților de intrare \mathbf{I} a unei rețele care minimizează eroarea E_c .

Evaluarea erorii din ecuația (4.91) indică faptul că aceasta conține termenii V_0^2 și V_1^2 , ceea ce face ca termenii w_{ij} să fie egali cu 2^{2i} în loc de zero. De aceea, se va adăuga un termen suplimentar de eroare E_a . În plus față de eliminarea elementelor diagonale ale matricii \mathbf{W} , acest termen va fi minimizat în apropierea colțurilor pătratului $[0, 1]$ din planul V_0, V_1 . De aceea, o alegere potrivită pentru E_a este:

$$E_a = -\frac{1}{2} \sum_{i=0}^1 2^{2i} V_i (V_i - 1) \quad (4.92)$$

Prin combinarea funcțiilor de eroare E_a și E_c rezultă următoarea funcție totală de energie, care este minimizată de convertorul A/D [190]:

$$E = \frac{1}{2} \left(x - \sum_{i=0}^1 V_i 2^i \right)^2 - \frac{1}{2} \sum_{i=0}^1 2^{2i} V_i (V_i - 1) \quad (4.93)$$

Primul termen reprezintă funcția de eroare a intrării x și a ieșirii binare (V_0, V_1). Al doilea termen este funcția de restricție care forțează ca ieșirea să fie 0 sau 1. Această expresie a energiei trebuie egalată cu formula generală a energiei, care este egală în acest caz cu:

$$E = -\frac{1}{2} \sum_{i=0}^1 \sum_{j \neq i, j=0}^1 w_{ij} V_i V_j - \sum_{i=0}^1 I_i V_i \quad (4.94)$$

Prin rearanjarea membrului drept al expresiei energiei din ecuația (4.94) se obține:

$$E = \frac{1}{2} x^2 + \frac{1}{2} \sum_{i=0}^1 \sum_{j \neq i, j=0}^1 2^{i+j} V_i V_j + \sum_{i=0}^1 (2^{2i-1} - 2^i x) V_i \quad (4.95)$$

Atunci când intrarea analogică x este fixată, termenul $(\frac{1}{2})x^2$ este o constantă aditivă în ecuația (4.95), și nu este relevantă pentru minimizarea funcției E indicate de ecuația (4.95) în planul V_0, V_1 ; de aceea, acest termen va fi omis. Comparând coeficienții asemănători ai V_0 și V_1 din partea dreaptă a funcțiilor de energie (4.94) și (4.95), rezultă următorii parametri ai rețelei:

$$\begin{aligned} w_{01} = w_{10} &= -2 \\ I_0 &= x - \frac{1}{2} \\ I_1 &= 2x - 2 \end{aligned} \quad (4.96)$$

Rețelele neuronale de tip gradient sunt exemple de sisteme neliniare, dinamice, și asimptotic stabile. Aceste rețele neuronale nu garantează însă convergența la minimumul global, din cauza faptului că nu dispun de posibilitatea unei căutări globale. Pentru a se evita blocarea în minime locale, unii autori au propus utilizarea procedurii de călire simulată pentru rețelele neuronale. Deoarece această procedură este secvențială, ea necesită un timp de calcul semnificativ pentru obținerea unei soluții optime. O altă problemă a acestor rețele este că un mare număr de parametri trebuie selectați și ajustați în mod corespunzător chiar și pentru a se asigura convergența la o soluție validă.

Avantajul rețelelor neuronale Hopfield este convergența rapidă la un minim stabil. Dacă la funcția de energie se adaugă condiții de restricție cu o pondere mare, punctul inițial va converge spre cel mai apropiat punct valid fără influență din partea funcției de cost. Aceasta diferă de cazul algoritmilor genetici, deoarece acești algoritmi nu pot ajusta configurația spre o soluție validă, chiar dacă penalizarea este severă. Această proprietate face ca rețelele neuronale Hopfield să poată fi utilizate pentru căutarea locală în cadrul algoritmilor genetici pentru rezolvarea problemelor combinatoriale și a problemelor cu restricții. Pe această bază, Lin *et al.* [114] au propus un model hibrid care utilizează algoritmi genetici și rețele neuronale Hopfield pentru optimizarea funcțiilor.

4.5.2.3 Utilizarea rețelelor neuronale pentru plasare

Considerăm cazul cel mai simplu al problemei de plasare. Fiind date n module de circuit și o matrice de conectivitate $\mathbf{C} = [C_{ij}]$, unde C_{ij} indică conectivitatea între modulul i și modulul j , cele n module trebuie plasate în n locații ale unei suprafețe bidimensionale, astfel încât lungimea totală de interconectare Manhattan să fie minimizată. Metoda se va ilustra prin circuitul din Figura 4.20(a), locațiile în care trebuie plasate modulele fiind ilustrate în Figura 4.20(b).

Soluția acestei probleme aparține lui Yu [187], care a utilizat rețelele neuronale Hopfield pentru rezolvarea problemei de plasare. Acesta a modificat soluția propusă de Hopfield și Tank pentru rezolvarea problemei comis-voiajorului. Se utilizează o rețea cu n^2 neuroni. Rețeaua constă dintr-o matrice de $n \times n$ neuroni, după cum se prezintă în Figura 4.21.

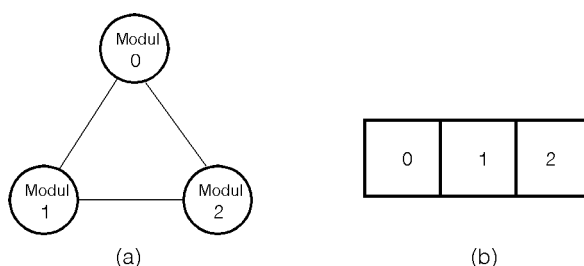


Figura 4.20. (a) Exemplu de graf al unui circuit. (b) Definiția pozițiilor.

Neuronii sunt numerotați de la 0 la $n^2 - 1$, de la stânga la dreapta și de sus în jos. Valoarea unui element din poziția (i, j) a matricii reprezintă "șansa" unui modul i de a fi poziționat în locația j . Fiecare linie corespunde unui modul de circuit. Fiecare coloană corespunde celor n locații posibile în care se pot plasa modulele. Pentru a obține o soluție fezabilă, un singur neuron din orice linie sau coloană poate avea ieșirea 1. Ieșirile neuronilor sunt normalizate, astfel încât acestea au valori între 0 și 1.

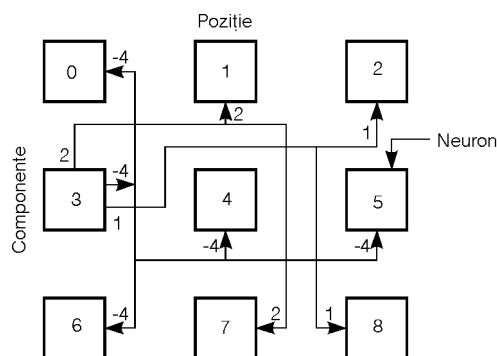


Figura 4.21. O rețea neuronală artificială utilizată pentru plasare.

Următoarea etapă este determinarea intensității sinapselor. Se calculează mai întâi distanța Manhattan între fiecare pereche de locații. Valoarea T_{k_1, j_1, j_2} a intensității sinapsei între neuronii k și l (care este un element al matricii intensității sinapselor) este definită ca și conectivitatea între modulele de circuit i_1 și i_2 multiplicată cu $f(j_1, j_2)$, unde f este o funcție a distanței între locațiile j_1 și j_2 , $k = i_1 \times \sqrt{n} + j_1$, iar $l = i_2 \times \sqrt{n} + j_2$. După experimentări funcția f s-a ales ca fiind egală cu (*offset* - distanța Manhattan dintre j_1 și j_2), unde parametrul *offset* este de obicei mai mare decât \sqrt{n} .

De exemplu, intensitatea sinapsei dintre neuronii 2 și 3 ($T_{2,3}$) se poate determina astfel. Neuronul 2 are $(i_1, j_1) = 0, 2$, iar neuronul 3 are $(i_2, j_2) = 1, 0$ (Figura 4.26). Astfel, $T_{2,3}$ va fi egal cu

$$C_{0,1} \times (\text{offset} - \text{distanța Manhattan dintre 2 și 0}) = 1 \times (3 - 2) = 1$$

unde s-a presupus *offset* = 3.

Conexiunile corespunzătoare ale neuronului 3 sunt prezentate în Figura 4.21.

Determinarea parametrilor rețelei neuronale (matricea ponderilor, valorile pragurilor, constantele implicate în funcția de energie și funcția de activare) este dificilă, soluția finală depinzând în mare măsură de acești parametri. De acești parametri depinde de asemenea convergența rețelei spre o soluție validă.

Rezultatele obținute de Yu în utilizarea rețelelor neuronale Hopfield pentru rezolvarea problemei de plasare nu au fost promițătoare. Unele din dificultățile indicate au fost timpii mari de simulare, calitatea redusă a soluțiilor și sensibilitatea soluțiilor de parametrii rețelei. Pe de altă parte, Sriram și Kang au obținut rezultate încurajatoare la plasarea bidimensională a unui număr de până la 64 de module, rețeaua utilizată fiind capabilă să determine soluțiile optime sau apropiate de cele optime în cele mai multe cazuri [190].

4.6 Algoritmi de plasare propuși pentru circuitele FPGA cu resurse limitate de rutare

O plasare adecvată pentru circuitele FPGA cu resurse limitate de rutare trebuie să realizeze nu numai minimizarea lungimii conexiunilor prin gruparea celulelor care sunt interconectate, dar trebuie să asigure și rutabilitatea circuitelor. În această secțiune se prezintă algoritmi de plasare elaborați pentru asemenea circuite FPGA, care au ca obiectiv asigurarea rutabilității circuitelor. În secțiunea 4.6.1 se prezintă un algoritm de plasare bazat pe metoda de partiționare prin tăietura minimă, care încearcă asigurarea rutabilității prin distribuirea uniformă a conexiunilor în cadrul partițiilor. În secțiunea 4.6.2 se prezintă un algoritm genetic pentru plasare, care alocă un număr de celule libere în procesul de plasare, care vor fi utilizate pentru rutare.

4.6.1 Algoritm de plasare pe baza tăieturii minime

4.6.1.1 Descrierea algoritmului de plasare

Tehnicile de plasare cele mai utilizate se bazează pe partiționarea prin metoda tăieturii minime. Acestea aplică în mod recursiv o procedură de bipartiționare cu scopul minimizării numărului de interconexiuni care intersectează linia de tăietură în fiecare etapă a algoritmului. Algoritmul realizează optimizarea prin plasarea apropiată a celulelor care sunt interconectate. De aceea, componenta principală a funcției de cost pentru un asemenea algoritm este lungimea interconexiunilor.

Se prezintă în continuare un algoritm de plasare bazat pe metoda de partiționare descrisă în capitolul 3 (secțiunea 3.7.1), cu o funcție obiectiv care urmărește pe lângă reducerea lungimii interconexiunilor, și distribuirea uniformă a conexiunilor în cadrul partițiilor. Ca urmare, rezultatul plasării maximizează posibilitatea unei rutări fezabile a circuitului. Este studiat de asemenea efectul utilizării diferitelor secvențe de aplicare a liniilor de tăietură, cu scopul reducerii numărului de conexiuni în apropierea centrului circuitului, deoarece în cazul plasării obișnuite care se bazează pe metoda tăieturii minime apare în mod frecvent creșterea numărului de conexiuni în această zonă.

Algoritmii de plasare care utilizează partiționarea ierarhică pe baza tăieturii minime au fost utilizați în mod frecvent. Fiind dat un circuit reprezentat printr-un graf, fiecare vârf reprezentând o celulă logică și fiecare muchie reprezentând o conexiune, acești algoritmi partiționează în mod recursiv circuitul și suprafața de plasare, până când se ajunge la un graf cu un singur nod, care va fi plasat într-o regiune a suprafeței. Obiectivul fiecărei etape de bipartiționare este minimizarea dimensiunii tăieturii, cu restricția ca fiecare porțiune să conțină același număr de noduri, sau un număr

apropiat. Această metodă are avantajul obținerii unei plasări de calitate corespunzătoare, a unui timp de execuție redus și a unei implementări simple.

Dezavantajul algoritmilor de plasare amintii este că dimensiunea tăieturii este singura metrică utilizată în cadrul funcției de cost. De aceea, este posibilă obținerea unei dimensiuni reduse a tăieturii, și în același timp a unor porțiuni cu un număr de conexiuni semnificativ diferit. Ca urmare, o asemenea plasare poate fi rutată într-un mod dificil, sau efectuarea rutării poate fi chiar imposibilă. Pentru a elimina acest dezavantaj, este necesar ca în cadrul funcției de cost să se ia în considerare nu numai dimensiunea tăieturii, ci și distribuția interconexiunilor din cele două porțiuni ale bipartiției.

Algoritmul de plasare propus utilizează bipartiționarea a cărei funcție obiectiv urmărește minimizarea lungimii interconexiunilor simultan cu minimizarea diferenței între numărul de conexiuni din cele două porțiuni. Această diferență, numită număr de dezechilibru, este măsura care urmărește distribuția echilibrată a conexiunilor din cele două porțiuni. În cadrul algoritmului de plasare se aplică în mod recursiv bipartiționarea cu funcția de cost descrisă în secțiunea 3.7.1, până când se ajunge la porțiuni care conțin o singură celulă a circuitului. Liniile de tăietură se aplică în mod alternativ, pe orizontală și pe verticală.

4.6.1.2 Secvența de aplicare a liniilor de tăietură

Pe lângă o procedură eficientă de partiționare, este necesară o strategie adecvată de aplicare a liniilor de tăietură. O secvență tradițională de aplicare a liniilor de tăietură, cum este procedura de partiționare quadratică [146], are dezavantajul că nu ține cont de poziția terminalelor externe. Pentru aceasta se poate utiliza tehnica numită propagarea terminalelor. Această tehnică are însă dezavantaje. De exemplu, cele două regiuni sunt procesate secvențial, neexistând criterii pentru stabilirea regiunii care trebuie procesată prima.

Secvența de aplicare a liniilor de tăietură are un rol important. În cadrul plasării convenționale bazată pe tăietura minimă, se alege linia de tăietură poziționată în centrul regiunii curente. Dacă liniile de tăietură se aplică în această ordine, Figura 4.22 indică secvența liniilor de tăietură, unde liniile întrerupte sunt cele aplicate recent, iar liniile continue sunt cele aplicate anterior.

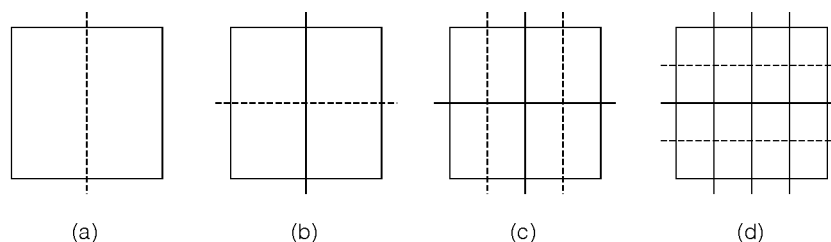


Figura 4.22. Alegerea liniei de tăietură poziționată în centru.

Presupunem că linia curentă de tăietură, cea indicată în Figura 4.23 printr-o linie întreruptă, este în imediată apropiere a liniei de tăietură din centru. Această linie va înjumătăți patru regiuni și conexiunile corespunzătoare. În fiecare etapă a acestei bipartiționări, se pot interschimba perechi de noduri, de exemplu astfel încât unul din noduri este în regiunea *A*, iar al doilea în *E*. Pentru a se ține cont de rezultatul bipartiționării pentru liniile de tăietură aplicate anterior, nu se permite interschimbarea nodurilor aflate în regiuni delimitate de liniile de tăietură aplicate anterior.

Numărul de noduri dintr-o regiune este proporțional cu suprafața regiunii respective. Deoarece linia de tăietură considerată este în apropierea liniei de tăietură din centru, regiunile intersectate au o suprafață redusă. De aceea, numărul nodurilor din aceste regiuni este redus. În consecință, numărul perechilor posibile care pot fi inter-schimbate în procesul de bipartiționare este limitat. De obicei, aceasta are ca rezultat o dimensiune a tăieturii relativ ridicată pentru liniile de tăietură din apropierea centrului, din cauza numărului mic de mutări posibile.

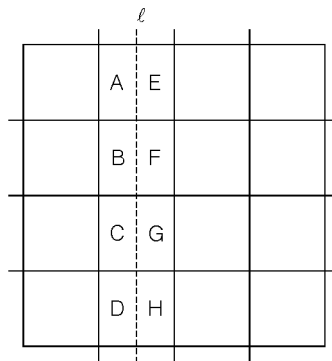


Figura 4.23. Plasare tradițională pe baza tăieturii minime.

Din cele de sus rezultă că, pentru a se reduce dimensiunea tăieturii în apropierea centrului, în această zonă liniile de tăietură trebuie aplicate în primele etape ale procesului de bipartiționare, după cum se indică în Figura 4.24. Pentru claritate, liniile de tăietură orizontale nu sunt indicate în această figură. Liniile de tăietură orizontale și verticale sunt aplicate alternativ și în acest caz.

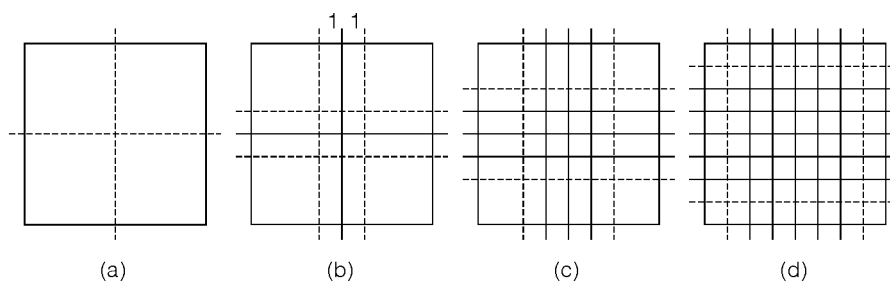


Figura 4.24. Secvența propusă pentru aplicarea liniilor de tăietură.

4.6.2 Algoritm genetic pentru plasarea circuitelor FPGA

Algoritmii genetici reprezintă o paradigmă eficientă pentru rezolvarea problemelor complexe de optimizare. După cum s-a arătat de J. H. Holland, simularea evoluției naturale în care populația trece printr-un proces de adaptare, cu o strategie controlată de evoluție, poate fi o alternativă eficientă pentru rezolvarea acestui tip de probleme. În timp ce alți algoritmi îmbunătățesc în mod iterativ o plasare inițială prin mutarea unei singure celule sau prin interschimbarea a două celule, un algoritm genetic pornește de la un set de plasări inițiale care reprezintă populația inițială. Algoritmul încearcă să combine caracteristicile adecvate din două configurații de plasare

pentru a forma o nouă plasare. Aceste caracteristici adecvate se numesc scheme, ele fiind plasările relative ale subseturilor de celule. Prin procesul de evoluție simulată, după un număr de generații, soluțiile candidate rețin caracteristicile mai bune ale soluțiilor multiple din generațiile anterioare. Acest paralelism intrinsec conferă un avantaj algoritmilor genetici față de metoda călirii simulate [102], care utilizează o singură soluție.

Datorită acestor avantaje, s-a încercat rezolvarea problemei de plasare a circuitelor FPGA printr-un algoritm genetic. După cunoștințele noastre, nu a fost publicat până în prezent un algoritm genetic pentru plasarea circuitelor FPGA. Algoritmul descris în această secțiune, implementat pentru seria de circuite *Atmel 6000*, reprezintă deci o primă contribuție în acest sens.

O plasare este acceptabilă dacă se poate realiza rutarea în proporție de 100%. Aceasta nu este o sarcină simplă pentru arhitecturile FPGA cu resurse limitate de rutare, cum sunt circuitele FPGA *Atmel 6000* [7]. O plasare corespunzătoare va realiza nu numai gruparea blocurilor conectate, dar va asigura ca blocurile logice să nu fie plasate foarte apropiat, pentru a se permite rutarea circuitului. O altă contribuție o reprezintă un algoritm care alocă un număr de celule libere în cadrul procesului de plasare, în scopul utilizării acestor celule pentru rutare. Funcția de cost utilizată optimizează un număr de diferite metrice, care cuprind atât lungimea interconexiunilor, cât și măsuri ale rutabilității plasării.

4.6.2.1 Utilizarea algoritmilor genetici pentru plasare

Deși nu au fost utilizați pentru plasarea circuitelor FPGA, algoritmii genetici au fost utilizați pentru plasarea circuitelor VLSI. Lucrarea clasică a fost realizată de Cohoon *et al.* [57] [58]. Aceștia au codificat o plasare prin notația poloneză a unui arbore binar, un cromozom fiind deci reprezentat printr-un șir. Au fost utilizați diferiți operatori de recombinare, care operează fie direct asupra șirurilor, fie iau în considerare structura de arbore prin decodificarea cromozomului.

Esbensen [73] a descris un algoritm genetic pentru plasarea macro-celulelor în care reprezentarea genotipului este de asemenea un arbore binar. Spre deosebire de abordarea lui Cohoon *et al.*, acest arbore nu caracterizează direct o plasare, ci aceasta poate fi generată prin decodificarea arborelui. Operatorii genetici lucrează direct cu structura arborelui. Esbensen și Mazumder [74] au raportat un algoritm numit SAGA, care este o generalizare a algoritmului genetic și a algoritmului de călire simulată. În funcție de setarea parametrilor săi de control, SAGA se comportă ca un algoritm genetic, un algoritm de călire simulată, sau o combinație a acestora. Autorii au arătat experimental că prin mixarea algoritmului genetic cu algoritmul de călire simulată se obține o plasare de calitate mai bună decât printr-un algoritm pur genetic.

Mohan și Mazumder [124] au descris un algoritm genetic distribuit pentru plasarea celulelor standard. Algoritmul a fost elaborat pentru a fi executat pe o rețea de stații de lucru. Procedura de plasare distribuită execută un algoritm genetic de bază pe fiecare procesor din rețea. Este introdus un nou operator genetic, *migrarea*, care transferă informații de plasare de la un procesor la altul în cadrul rețelei.

Schnecke și Vornberger [150] au prezentat un algoritm genetic pentru proiectarea fizică a circuitelor VLSI. Algoritmul combină etapele de amplasare și de rutare, optimizând simultan plasarea celulelor și rutarea. Aceiași autori au descris în [149] un algoritm genetic paralel pentru optimizarea combinată a plasării și a rutării. Procesul de căutare este auto-adaptabil. Se execută mai mulți algoritmi genetici secvențiali, cu strategii diferite. La intervale fixe de timp, aceste strategii sunt evaluate și fiecare strategie este ajustată pentru a obține rezultate mai bune. Ambii algoritmi [149], [150] utilizează un genotip codificat ca un arbore binar, care definește plasarea relativă a celulelor.

Considerăm graful din Figura 4.25(a) ale cărui vârfuri reprezintă module, iar numerele din dreptul muchiilor reprezintă ponderi ale interconexiunilor. Cele nouă module pot fi plasate în nouă locații, care sunt definite în Figura 4.25(b). O soluție posibilă a problemei de plasare este indicată în Figura 4.25(c).

Soluția din Figura 4.25(c) poate fi reprezentată ca un șir de simboluri. Dacă primul simbol din stânga șirului corespunde locației 0 din Figura 4.25(b), iar simbolul din dreapta șirului corespunde locației 8, soluția din Figura 4.25(c) poate fi reprezentată prin șirul [AGHCBIDEF] $\left(\frac{1}{85}\right)$. Numărul din paranteză reprezintă valoarea funcției de viabilitate, care este inversa lungimii ponderate a conexiunilor pe baza metricii Manhattan [146].

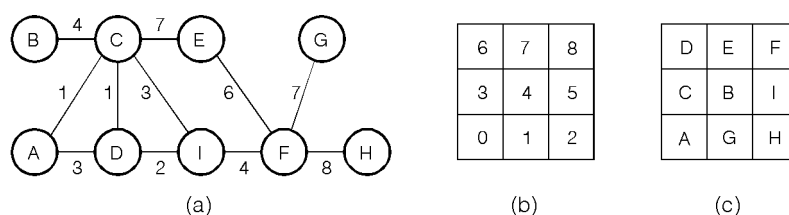


Figura 4.25. (a) Graful unui circuit. (b) Definiția pozițiilor. (c) O plasare posibilă.

Dacă locația din colțul din stânga jos al Figurii 4.25(b) este considerată în originea sistemului de coordonate cartezian, se pot calcula în mod simplu coordonatele oricărui modul pe baza indexului în șirul care reprezintă soluția plasării. De exemplu, indexul modulului E este 7. Coordonatele carteziene ale acestui modul sunt date de $x = (7 \bmod 3) = 1$, iar $y = \left\lceil \frac{7}{3} \right\rceil = 2$.

Orice șir de lungime 9 care conține caracterele [A, B, C, D, E, F, G, H, I] reprezintă o soluție posibilă. Alte soluții posibile (cromozomi) sunt [BDEFIGCHA] $\left(\frac{1}{110}\right)$, [IHAGBFCED] $\left(\frac{1}{95}\right)$, [BIDEFAGHC] $\left(\frac{1}{86}\right)$.

Și în cazul plasării, operatorii genetici cei mai utilizați sunt încrucișarea, mutația și inversiunea.

Cel mai simplu operator de *încrucișare* constă în alegerea aleatoare a unui punct de tăietură și generarea urmașului prin combinarea segmentului din stânga punctului de tăietură al unui părinte cu segmentul din dreapta punctului de tăietură al celuilalt părinte. Din exemplul anterior, considerăm cei doi părinți [BIDEF.AGHC] $\left(\frac{1}{86}\right)$ și [BDEFI.GCHA] $\left(\frac{1}{110}\right)$. Dacă punctul de tăietură este ales după poziția 4, urmașul produs va fi [BIDEF.GCHA]. Lungimea ponderată a conexiunilor pentru urmașul generat este redusă la 63, și astfel funcția de viabilitate a urmașului are valoarea $\frac{1}{63}$.

În exemplul anterior, elementele de la stânga punctului de tăietură al unui părinte nu apar la dreapta punctului de tăietură al celuilalt părinte. În caz contrar, unele simboluri din șirul soluției ar fi repetate. În cazul plasării celulelor aceasta nu ar reprezenta o soluție fezabilă. Există diferite modificări ale operației de încrucișare descrise anterior în scopul evitării repetiției simbolurilor, de exemplu: (a) încrucișare ordonată, (b) încrucișare mapată parțial (PMX), și (c) încrucișare ciclică.

Se exemplifică două operații de încrucișare utilizate în sistemul de plasare *Genie*, care a fost elaborat pentru plasarea modulelor într-o rețea rectangulară [57]. Primul operator de încrucișare selectează un modul aleator e_s din părintele 1 și mută acest modul și vecinii săi din acest părinte în locațiile corespunzătoare vecinilor modulului din părintele 2. Modulele care ocupau locațiile vecinilor din părintele 2 sunt deplasate apoi cu câte o locație în direcția vechilor locații ale modulelor care au fost

mutate până când este găsită o locație liberă. Aceste deplasări sunt ilustrate în Figura 4.26. Rezultatul este că se copiază din părintele 1 în părintele 2 o porțiune conținând e_s și cei patru vecini ai săi, și că alte module sunt deplasate cu cel mult o poziție.

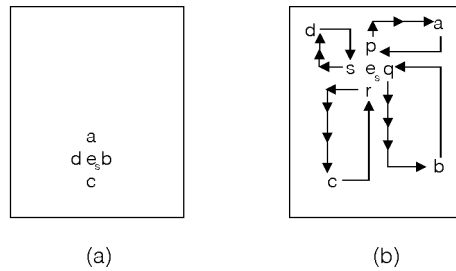


Figura 4.26. (a) Un modul aleator și vecinii acestuia. (b) Vecinii din (a) ai părintelui 1 înlocuiesc modulele vecine ale părintelui 2.

Al doilea operator de încrucișare utilizat în sistemul *Genie* selectează un pătrat conținând $k \times k$ module din părintele 1 și îl copiază în părintele 2. Numărul k este ales în mod aleator, având o medie de 3. Această metodă are tendința să dubleze unele module. De exemplu, în Figura 4.27, dacă modulele din pătratul părintelui 1 sunt copiate în pătratul părintelui 2, acele module din pătratul părintelui 1 care nu se află în pătratul părintelui 2 vor fi duplicate. Această problemă este rezolvată astfel: fie $SP_2 - SP_1$ setul modulelor din pătratul părintelui 2 care nu se află în pătratul părintelui 1 ($SP_2 - SP_1 = \{x, w, p, m\}$). Similar, $SP_1 - SP_2$ este setul modulelor din pătratul părintelui 1 care nu se află în pătratul părintelui 2 ($SP_1 - SP_2 = \{c, d, e, f\}$). Fiecare modul din $SP_2 - SP_1$ este mutat într-o locație ocupată în acel moment de un modul din $SP_1 - SP_2$. Deci, în părintele 2 modulele $\{x, w, p, m\}$ sunt mutate în locațiile ocupate de modulele $\{c, d, e, f\}$. Toate modulele din pătratul părintelui 1 sunt copiate apoi în pătratul părintelui 2 pentru a se obține un nou urmaș.

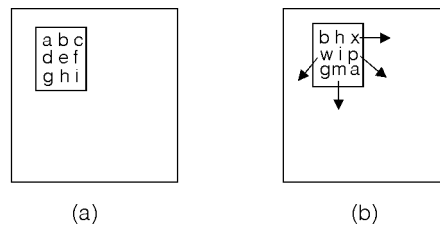


Figura 4.27. (a) Este selectat un pătrat în părintele 1. (b) Modulele pătratului din părintele 1 sunt copiate în părintele 2 și modulele duplicate sunt eliminate.

Mutația produce modificări incrementale aleatoare ale urmașului generat prin încrucișare. În cazul plasării, mecanismul de mutație cel mai utilizat constă în interschimbarea perechilor. Pentru plasare, o genă constând dintr-un triplet ordonat format dintr-o celulă și coordonatele asociate acesteia nu poate fi prezentă în nici un alt individ al populației. În acest caz, nu este utilă utilizarea doar a operației de încrucișare, deoarece aceasta reprezintă un mecanism de moștenire. Operatorul de mutație generează noi tripleți celulă-coordonate. Dacă noii tripleți sunt corespunzători, configurațiile care le conțin sunt reținute. Mutația este controlată de un parametru numit rata mutației M_r . O rată redusă a mutației înseamnă că infuzia de noi gene este foarte scăzută. O rată M_r ridicată va determina ca urmașii să piardă asemănarea cu părinții, ceea ce înseamnă că algoritmul se comportă ca un proces fără memorie, pierzându-și abilitatea de a învăța din istoria căutării [153].

În cazul operației de *inversiune*, se aleg în mod aleator două puncte de tăietură în cadrul cromozomului, și se inversează simbolurile subșirului din secțiunea tăiată. De exemplu, șirul [BID.EFGCH.A] (tăiat după pozițiile cu indexul 2 și 7) va deveni după inversiune [BID.HCGFE.A]. Operația este executată astfel încât nu modifică soluția reprezentată de cromozom, ci modifică numai reprezentarea acestuia. Deci, simbolurile șirului trebuie să aibă o interpretare independentă de poziția lor [83] [153].

Până acum s-au considerat operații asupra unor soluții posibile complete. O soluție parțială este reprezentată de o schemă. Pentru exemplul din Figura 4.25, șirul [*I*EFH*G*] indică o plasare parțială, celulele din pozițiile 0, 2, 6, 8 fiind nedefinite. Deoarece celula F este puternic conectată cu celulele I, E, G și H, iar aceste celule sunt adiacente cu celula F în soluția parțială amintită, această schemă reprezintă o plasare parțială de calitate. O asemenea plasare parțială va îmbunătăți funcția de viabilitate a unui individ care o moștenește. Astfel, o schemă este o plasare topologică a celulelor corespunzătoare, indicând pozițiile lor relative și nu locațiile lor exacte.

4.6.2.2 Reprezentarea soluției

Reprezentarea soluției are un rol major în proiectarea unui algoritm genetic eficient. Principalele caracteristici ale unei soluții sunt următoarele [41]:

- Este întotdeauna un șir de gene;
- Trebuie să permită o mare diversitate într-o populație de dimensiuni reduse;
- Trebuie să permită aplicarea simplă a operatorilor genetici;
- Trebuie să permită calcularea simplă a funcției obiectiv.

Pentru soluția problemei de plasare s-a ales reprezentarea printr-un șir de înregistrări. Fiecare înregistrare reprezintă o celulă, conținând identificatorul celulei și coordonatele (x, y) ale acesteia.

Poziția înregistrării corespunzătoare unei celule nu determină întotdeauna poziția fizică a celulei în circuit. Totuși, în anumite puncte din algoritm, poziția celulei este recalculată pe baza ordonării înregistrărilor celulelor în șir, după cum urmează. Pornind de la prima linie de celule din circuit, celulele sunt listate în ordine de la stânga la dreapta. La sfârșitul liniei, se trece la linia următoare, iar celulele sunt listate în ordine inversă. Procesul este continuat până când sunt listate toate celulele, schimbându-se direcția după fiecare linie. Fiind dat un șir de înregistrări, celulelor li se pot asigna poziții în cadrul liniilor, prin inversarea procesului anterior.

4.6.2.3 Operatori genetici

Încrușișarea este principalul operator genetic. Dacă Π_a și Π_b sunt doi indivizi părinți, reprezentând plasări valide, operația de încrușișare este definită ca $\Pi_a \times \Pi_b \rightarrow \Pi_u$, unde Π_u reprezintă urmașul, fiind o altă plasare validă. Încrușișarea combină scheme din ambii părinți, și astfel urmașul moștenește unele din caracteristicile părinților.

Cea mai simplă formă de încrușișare, descrisă în secțiunea 4.6.2.1, nu poate fi aplicată în acest caz, deoarece poate crea șiruri care nu au un corespondent fizic. De exemplu, o încrușișare simplă între două șiruri ABCD și BDCA poate produce două noi șiruri ABCA și BDCD. Fiecare din acestea are două caractere care se repetă și un caracter lipsă. Dacă fiecare caracter corespunde unei celule în cadrul problemei de plasare, atunci într-un șir valid fiecare caracter trebuie să apară o singură dată. În consecință, sunt necesari operatori de încrușișare mai complecși, care păstrează corectitudinea plasării. În algoritmul descris, s-a utilizat un operator numit încrușișare mapată parțial (*PMX*).

Încrucișarea *PMX* se execută astfel: se selectează doi părinți (1 și 2) și se alege un punct de tăietură arbitrar. Ca și în cazul încrucișării simple, întregul subșir din dreapta al părintelui 2 este copiat în șirul urmașului. Apoi, subșirul din stânga al părintelui 1 este parcurs genă cu genă, de la stânga până la punctul de tăietură. Dacă o genă nu există în cadrul urmașului, este copiată în șirul acestuia. Dacă însă gena există deja în cadrul urmașului, este determinată poziția sa în părintele 2, și este copiată gena din poziția determinată a părintelui 1.

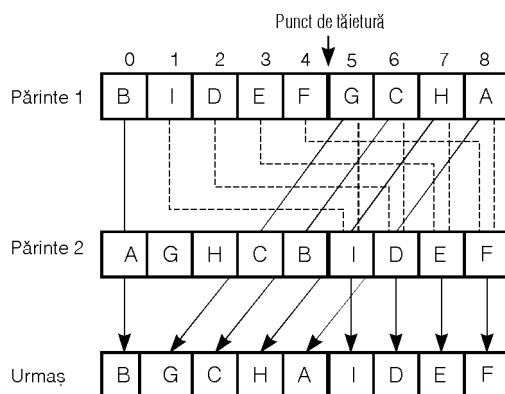


Figura 4.28. Încrucișarea *PMX*.

Figura 4.28 exemplifică modul în care se execută încrucișarea *PMX*. Fie punctul de tăietură după poziția 4. Subșirul din dreapta al părintelui 2, care este IDEF, este copiat în cadrul urmașului. Apoi, primul părinte este parcurs de la stânga, și deoarece gena B (poziția 0) nu există în cadrul urmașului, este copiată în poziția 0. Următoarea genă, I (poziția 1) există în poziția 5. Gena din poziția 5 a părintelui 1 este G, care nu există în cadrul urmașului, și de aceea gena G este copiată în poziția 1. Parcurgerea primului părinte se continuă în același mod, urmașul generat fiind BGCHAIDEF.

Este esențial ca indivizii cu viabilitate mai ridicată să participe la reproducere mai frecvent. În cazul în care un individ se reproduce de un număr mare de ori în primele etape ale algoritmului genetic, este probabil că toate soluțiile din populație vor moșteni anumite trăsături ale acestui individ și astfel vor fi asemănători. De aceea, strategia de selecție a părinților trebuie să evite această convergență prematură a algoritmului genetic la un optim local. În cadrul algoritmului aceasta se obține prin faptul că unui părinte cu viabilitate ridicată nu i se permite reproducerea decât o singură dată într-o generație.

Mutația produce modificări aleatoare incrementale ale urmașului generat prin încrucișare. Pentru problema de plasare, ca tehnică de mutație se poate utiliza interschimbarea perechilor. Este posibil ca o genă constând dintr-un triplet ordonat al unei celule și coordonatele ideale asociate acesteia să nu fie prezentă în niciunul din indivizii populației. În acest caz, încrucișarea nu este suficientă, deoarece este doar un mecanism de moștenire. Operatorul de mutație generează noi tripleți celulă-coordonate. În cazul în care aceștia se comportă în mod corespunzător, configurațiile care le conțin vor fi reținute.

Mutația este controlată de un parametru numit rata de mutație M_r . O valoare redusă a acestei rate înseamnă că infuzia noilor gene este foarte redusă. Un asemenea nivel redus de mutație va modifica structura soluției numai într-o mică măsură. O valoare ridicată a ratei de mutație va determina ca urmașii să-și piardă asemănarea cu părinții acestora. Procesul de mutație va permite ca diversitatea populației să fie menținută în etapele finale ale algoritmului genetic. Mutația permite de asemenea ca algoritmul genetic să evite minimele locale.

Inversiunea este o operație care modifică lungimea efectivă a unei scheme fără a altera viabilitatea individului în scopul creșterii probabilității de supraviețuire a schemelor mai lungi. Operatorul de inversiune modifică pozițiile înregistrărilor celulelor în șirul care reprezintă o plasare, dar nu modifică poziția fizică a celulelor din circuit. Inversiunea este controlată de un parametru numit rata de inversiune I_r .

4.6.2.4 Selecția populației pentru generația următoare

Algoritmii genetici convenționali utilizează metoda prin care se generează doi urmași din doi părinți și se înlocuiesc părinții prin cei doi urmași pentru prelucrările din generația următoare. În algoritmul propus s-a utilizat însă o strategie diferită. După generarea urmașilor, aceștia înlocuiesc un număr echivalent de indivizi din populația curentă, aleși dintre cei cu viabilitatea cea mai redusă. Această strategie permite supraviețuirea soluțiilor mai bune de-a lungul unui mare număr de generații. Soluțiile generate prin operatorii genetici concurează pentru cel puțin o generație. Această metodă de selecție a populației următoare asigură performanțe corespunzătoare ale algoritmului genetic.

4.6.2.5 Aspecte specifice pentru circuitele FPGA cu resurse limitate de rutare

Pentru circuitele FPGA cu resurse limitate de rutare, încă în etapa de plasare sunt necesare măsuri speciale pentru a se asigura rutabilitatea circuitelor. Funcția de cost utilizată optimizează diferite metrici, care includ lungimea totală a interconexiunilor, ca și măsuri ale rutabilității plasării. Pentru calculul lungimii interconexiunilor s-a utilizat metoda semi-perimetrului. În general, minimizarea lungimii interconexiunilor are ca efect plasarea apropiată a celulelor în centrul circuitului, ceea ce va conduce aproape sigur la conexiuni care nu pot fi rutate. Pentru a se asigura rutabilitatea plasării, ideea este de a se utiliza un număr de celule libere pentru rutare. Aceste celule sunt alocate în cadrul procesului de plasare.

În timp ce minimizarea lungimii totale a interconexiunilor reduce numărul resurselor de rutare necesare în mod global, nu poate asigura faptul că semnalele vor putea fi rutate local, datorită resurselor limitate de rutare. Pentru rezolvarea acestei probleme, au fost adăugate două componente la funcția de cost. Componenta de rutabilitate locală atribuie o penalizare acelor situații în care se poate determina faptul că o conexiune nu poate fi rutată utilizând resurse locale de rutare. Fiecare funcție necesită două sau trei intrări, din care numai una poate fi furnizată printr-o magistrală locală sau expres [7]. De aceea, funcțiile cu două intrări trebuie să recepționeze una din intrări de la celulele vecine, iar funcțiile cu trei intrări trebuie să recepționeze două din intrări de la aceste celule. Există patru celule adiacente care pot furniza aceste intrări, iar funcția de rutabilitate locală testează dacă semnalele de intrare necesare sunt fie prezente în aceste celule, fie există suficiente resurse de rutare disponibile astfel încât ele pot fi rutate. Rutabilitatea locală poate depista numai unele plasări ilegale care pot fi determinate din contextul imediat.

Componenta de echilibrare a densității este prevăzută pentru a se preveni congestia rutării datorită unei concentrări ridicate a funcțiilor într-o zonă a circuitului. Aceasta se realizează prin considerarea unor ferestre de celule (de 3×3 celule) și controlizarea intrărilor utilizate în această regiune. Pentru a se asigura distribuția uniformă a celulelor libere, penalizarea utilizată este pătratul numărului de intrări utilizate peste un anumit prag într-o fereastră, valori care se însumează pentru toate ferestrele. Ridicarea la pătrat este necesară pentru a se penaliza acele situații în care intrările utilizate au o concentrație ridicată. Pragul este necesar pentru ca o logică de dimensiuni reduse să nu fie răspândită în cadrul circuitului.

Se examinează fiecare fereastră din circuitul FPGA, astfel încât ferestrele se suprapun. Este permisă deplasarea ferestrelor dincolo de limitele circuitului, pentru ce-

lulele virtuale aflate dincolo de marginile circuitului presupunându-se un număr de intrări utilizate egal cu media globală. În cazul în care nu s-ar permite deplasarea ferestrelor dincolo de limitele circuitului, sau s-ar presupune că celulele virtuale nu au intrări utilizate, un mare număr de celule ar fi concentrate la marginile circuitului. Similar, dacă s-ar presupune că celulele virtuale ar avea toate intrările utilizate, ar fi evitată alocarea celulelor de la margine.

Descrierea formală a algoritmului genetic pentru plasare este prezentată în Figura 4.29.

```

Algorithm Plasare_Genetică;
  /*  $N_p$  - dimensiunea populației */
  /*  $N_g$  - numărul de generații */
  /*  $N_u$  - numărul de urmași */
  /*  $C_r$  - rata de încrucișare */
  /*  $I_r$  - rata de inversiune */
  /*  $M_r$  - rata de mutație */

begin
  Se citește lista de conexiuni a circuitului;
  Se generează o populație inițială de  $N_p$  configurații de plasare;
  for  $j = 1$  to  $N_p$  do
    Se evaluează Fitness (Populație [ $j$ ]);
  endfor;
   $N_u = C_r * N_p$ ;
  for ( $i = 1$  to  $N_g$ ) do
    for ( $j = 1$  to  $N_u$ ) do
      /* Se selectează părinții pentru reproducție cu o probabilitate
        proporțională cu valoarea viabilității */
       $(x, y) \leftarrow Sel\_părinți$ ;
      /* Se aplică operatorul de încrucișare */
       $urmaș [j] \leftarrow Generare\_urmaș (x, y)$ ;
      Se aplică cu rata  $M_r$  operatorul de mutație asupra
      configurațiilor selectate;
      Se aplică cu rata  $I_r$  operatorul de inversiune asupra
      configurațiilor selectate;
    endfor;
    for  $j = 1$  to  $N_p$  do
      Se evaluează Fitness (Populație [ $j$ ]);
    endfor;
    Se determină viabilitatea medie a populației;
    Se selectează  $N_p$  configurații pentru următoarea generație, configurațiile
    cu viabilitate ridicată având o probabilitate mai mare de selecție;
  endfor;
  Se returnează configurația cu viabilitatea cea mai ridicată;
end.

```

Figura 4.29. Algoritm genetic pentru plasare.

Timpul de execuție τ al algoritmului este o funcție de dimensiunea populației N_p , numărul de generații N_g , numărul de celule n , și valorile asignate parametrilor genetici.

$$\tau(n, N_p, N_g) = t_{init}(n, N_p) + N_g \cdot (t_{rep}(n, N_p) + t_{eval}(n, N_p)) + t_{out}(n) \quad (4.97)$$

unde

t_{init} timpul necesar pentru citirea listei de conexiuni a circuitului și pentru generarea populației inițiale;

- t_{rep} timpul necesar pentru execuția tuturor operațiilor genetice (pentru o generație) și pentru crearea unui set de noi indivizi;
- t_{eval} timpul necesar pentru evaluarea funcției de cost pentru fiecare nou individ și selecția indivizilor care vor supraviețui în următoarea generație;
- t_{out} timpul necesar pentru scrierea soluției finale pe disc.

Timpul cu ponderea cea mai mare este t_{eval} , iar $N_g \cdot t_{eval}$ este în jur de 75% din timpul total τ .

4.6.3 Rezultate experimentale

Algoritmul de plasare propus care utilizează secvența propusă a liniilor de tăietură a fost implementat în limbajul C. Experimentele au fost efectuate pe un calculator IBM PC cu un procesor Pentium de 133 MHz, sub sistemul de operare Windows NT Version 4.0. S-a utilizat un număr de nouă circuite de test din cadrul setului de circuite al centrului MCNC (*Microelectronics Center of North Carolina*).

Experimentele au fost efectuate după cum urmează. Lista de conexiuni a circuitelor de test a fost convertită din formatul EDIF (*Electronic Design Interchange Format*) în formatul utilizat de programul de mapare tehnologică. Circuitele au fost apoi mapate tehnologic pentru circuitul FPGA *Atmel 6002*. Algoritmul de plasare a fost aplicat asupra listelor de conexiuni obținute după maparea tehnologică.

În Tabelul 4.1 s-a comparat suma dimensiunilor tăieturilor pentru toate liniile de tăietură aplicate pentru algoritmul propus și algoritmul tradițional. Această sumă este o măsură a lungimii totale a conexiunilor atunci când estimarea se realizează prin metoda semiperimetrului. Compararea s-a efectuat pentru două cazuri. Primul caz este cel în care nu se urmărește echilibrarea numărului de conexiuni din cele două părți ($PONDERE = 0$), caz indicat în tabel prin $P = 0$. Al doilea caz este cel în care se realizează și echilibrarea numărului de conexiuni, termenul corespunzător din funcția de cost având aceeași pondere cu termenul pentru minimizarea lungimii conexiunilor ($PONDERE = 1$), caz indicat în tabel prin $P = 1$. Valorile s-au obținut prin rularea algoritmilor de 20 de ori, calculându-se media sumei tăieturilor pentru aceste rulări (valorile fiind rotunjite la numere întregi).

Tabelul 4.1. Media sumei tăieturilor pentru algoritmul de plasare propus și pentru algoritmul tradițional.

Circuit	P = 0			P = 1		
	Alg. propus	Alg. tradițional	Reducere	Alg. propus	Alg. tradițional	Reducere
<i>b1</i>	13	19	31.5 %	15	17	11.7 %
<i>c17</i>	8	10	20.0 %	8	10	20.0 %
<i>cm138a</i>	23	38	39.4 %	28	38	26.3 %
<i>con1</i>	17	30	43.3 %	17	28	39.2 %
<i>decod</i>	44	69	36.2 %	51	59	13.5 %
<i>majority</i>	11	17	35.2 %	11	15	26.6 %
<i>tcon</i>	45	60	25.0 %	41	51	19.6 %
<i>x2</i>	38	51	25.4 %	42	46	8.6 %

Din Tabelul 4.1 rezultă că prin algoritmul propus s-a obținut o reducere a sumei tăieturilor cu un procent cuprins între 20 % și 43.3 % pentru cazul în care nu se urmărește echilibrarea numărului de conexiuni ($P = 0$). În acest caz, reducerea medie este de 32 %. Pentru cazul în care se urmărește și echilibrarea numărului de conexiuni ($P = 1$), reducerea obținută este cuprinsă între 8.6 % și 39.2 %, reducerea medie fiind de 20.6 %. Faptul că reducerea obținută este mai mare în cazul în care nu se urmărește echilibrarea numărului de conexiuni este justificat, deoarece singura metrică în acest caz este dimensiunea tăieturii.

În Tabelul 4.2 s-a comparat valoarea maximă a tăieturii pentru toate liniile de tăietură aplicate pentru algoritmul propus și algoritmul tradițional. Compararea s-a efectuat pentru aceleași două cazuri: $P = 0$ și $P = 1$.

Tabelul 4.2. Valorile maxime ale tăieturilor pentru algoritmul de plasare propus și pentru algoritmul tradițional.

Circuit	P = 0			P = 1		
	Alg. propus	Alg. tradițional	Reducere	Alg. propus	Alg. tradițional	Reducere
<i>bl</i>	7	9	22.2 %	8	11	27.2 %
<i>c17</i>	4	4	0.0 %	4	6	33.3 %
<i>cm138a</i>	9	12	25.0 %	14	16	12.5 %
<i>con1</i>	9	12	25.0 %	8	10	20.0 %
<i>decod</i>	10	12	16.6 %	15	17	11.7 %
<i>majority</i>	7	10	30.0 %	9	14	35.7 %
<i>tcon</i>	10	11	9.0 %	14	17	17.4 %
<i>x2</i>	12	14	14.2 %	12	16	25.0 %

Conform Tabelului 4.2, prin algoritmul propus s-a obținut o reducere a tăieturii maxime cu un procent de până la 30 % pentru cazul în care nu se urmărește echilibrarea numărului de conexiuni. În acest caz, reducerea medie este de 17.75 %. Pentru cazul în care se urmărește și echilibrarea numărului de conexiuni, reducerea obținută este de până la 35.7 %, reducerea medie fiind de 22.85 %. Deci, prin algoritmul propus se obține atât o reducere a sumei tăieturilor, cât și o reducere a valorii maxime a tăieturii.

Algoritmul genetic pentru plasare a fost implementat în limbajul C. Experimentele au fost efectuate utilizând de asemenea circuite de test din setul *MCNC*. S-a studiat efectul dimensiunii populației și a ratei de mutație asupra calității rezultatelor în scopul determinării unor valori optime pentru acești parametri. Calitatea rezultatelor a fost apreciată pe baza lungimii totale a interconexiunilor. Experimentele au arătat că nu se obțin diferențe semnificative a lungimii totale a interconexiunilor cu variația ratei de mutație, dar există diferențe în funcție de dimensiunea populației. Tabelele 4.3 - 4.11 indică lungimea totală a interconexiunilor în starea inițială și cea finală în funcție de dimensiunea populației (rata de mutație fiind constantă, $M_r = 0.05$), și timpul de execuție al algoritmului pentru un număr de generații $N_g = 1000$.

Tabelul 4.3. Efectul dimensiunii populației asupra lungimii totale a interconexiunilor pentru circuitul *b1*.

Dimens. populație	Lungimetotală a conexiunilor		Timp de execuție (s)
	Inițială	Finală	
20	770	425	11
40	742	379	22
60	726	300	37
80	726	267	55
100	726	268	78
120	723	310	106

Tabelul 4.4. Efectul dimensiunii populației asupra lungimii totale a interconexiunilor pentru circuitul *c17*.

Dimens. populație	Lungime totală a conexiunilor		Timp de execuție (s)
	Inițială	Finală	
20	363	184	7
40	363	163	14
60	363	129	22
80	363	105	32
100	355	149	44
120	355	129	56

Tabelul 4.5. Efectul dimensiunii populației asupra lungimii totale a interconexiunilor pentru circuitul *cm138a*.

Dimens. populație	Lungime totală a conexiunilor		Timp de execuție (s)
	Inițială	Finală	
20	1201	677	17
40	1154	738	38
60	1154	543	59
80	1154	554	88
100	1154	519	120
120	1154	571	160

Tabelul 4.6. Efectul dimensiunii populației asupra lungimii totale a interconexiunilor pentru circuitul *con1*.

Dimens. populație	Lungime totală a conexiunilor		Timp de execuție (s)
	Inițială	Finală	
20	899	458	12
40	899	498	26
60	899	388	43
80	899	330	61
100	887	387	84
120	887	318	109

Tabelul 4.7. Efectul dimensiunii populației asupra lungimii totale a interconexiunilor pentru circuitul *daio*.

Dimens. populație	Lungime totală a conexiunilor		Timp de execuție (s)
	Inițială	Finală	
20	901	450	11
40	734	343	24
60	734	399	39
80	734	293	56
100	734	334	78
120	734	401	102

Tabelul 4.8. Efectul dimensiunii populației asupra lungimii totale a interconexiunilor pentru circuitul *decod*.

Dimens. populație	Lungime totală a conexiunilor		Timp de execuție (s)
	Inițială	Finală	
20	2485	1364	34
40	2485	1301	80
60	2485	1338	134
80	2485	1289	199
100	2485	1317	279
120	2485	1497	376

Tabelul 4.9. Efectul dimensiunii populației asupra lungimii totale a interconexiunilor pentru circuitul *majority*.

Dimens. populație	Lungime totală a conexiunilor		Timp de execuție (s)
	Inițială	Finală	
20	494	250	8
40	339	219	16
60	339	161	26
80	339	140	38
100	339	142	52
120	339	121	68

Tabelul 4.10. Efectul dimensiunii populației asupra lungimii totale a interconexiunilor pentru circuitul *tcon*.

Dimens. populație	Lungime totală a conexiunilor		Timp de execuție (s)
	Inițială	Finală	
20	2267	1388	33
40	2264	1215	77
60	2264	1242	126
80	2264	1296	185
100	2264	1249	255
120	2172	1323	338

Tabelul 4.11. Efectul dimensiunii populației asupra lungimii totale a interconexiunilor pentru circuitul *x2*.

Dimens. populație	Lungime totală a conexiunilor		Timp de execuție (s)
	Inițială	Finală	
20	2258	1264	28
40	2201	1240	63
60	2201	1162	103
80	2192	1211	154
100	2192	1090	213
120	2192	1246	285

Din aceste tabele rezultă că cele mai bune valori s-au obținut pentru dimensiuni ale populației de 80 sau 100. Pe baza acestor experimente, au fost utilizate următoarele valori ale parametrilor algoritmului genetic: dimensiunea populației $N_p = 80$, rata de mutație $M_r = 0.05$, rata de inversiune $I_r = 0.15$, numărul de generații $N_g = 1000$. Rezultatele obținute pentru acești parametri sunt prezentate în Tabelul 4.12.

Tabelul 4.12. Lungimea totală a interconexiunilor și timpul de execuție pentru setul de circuite de test.

Circuit	Lungime totală a conexiunilor		Timp de execuție (s)
	Inițială	Finală	
<i>bl</i>	726	267	55
<i>c17</i>	363	105	32
<i>cm138a</i>	1154	554	88
<i>con1</i>	889	330	61
<i>daio</i>	734	293	56
<i>decod</i>	2485	1289	199
<i>majority</i>	339	140	38
<i>tcon</i>	2264	1296	185
<i>x2</i>	2192	1211	154

4.7 Concluzii

În acest capitol a fost studiată problema de plasare a modulelor, având ca principal obiectiv asigurarea rutabilității circuitelor. Au fost prezentate principalele funcții de cost și restricții utilizate pentru rezolvarea acestei probleme. Funcția obiectiv cea mai des utilizată în cadrul plasării este lungimea totală estimată a conexiunilor. Au fost prezentate diferite tehnici pentru estimarea lungimii totale a interconexiunilor pentru o plasare dată. Au fost descrise și alte funcții de cost utilizate, ca minimizarea tăieturii maxime și minimizarea densității maxime.

Pentru rezolvarea problemei de plasare s-a utilizat un număr mare de tehnici euristice. Metodele bazate pe partiționare implică aplicarea recursivă a unui algoritm de partiționare, de obicei algoritmul Kernighan-Lin sau algoritmul de călire simulată. Deși metoda de călire simulată poate obține soluția optimă globală, timpul de calcul necesar pentru circuite de dimensiuni mari este foarte mare. Pentru a elimina dezavantajul complexității ridicate a algoritmului de călire simulată, au fost propuse mai multe tehnici de creștere a vitezei pentru acest algoritm, dintre care au fost descrise cele propuse de Lam și Delosme, respectiv Mallela și Grover.

O altă metodă care a fost descrisă utilizează partiționarea ierarhică pe baza tăieturii proporționale. Circuitul este descompus mai întâi într-un arbore de grupuri prin aplicarea recursivă a metodei de partiționare pe baza tăieturii proporționale. În etapa a doua, fiecare grup este considerat ca o super-componentă, și se aplică metoda de călire simulată pentru acestea. În etapa a treia, se integrează soluțiile subproblemelor prin utilizarea metodei ferestrelor mobile. Prin această metodă s-a obținut o anumită reducere a dimensiunii maxime a tăieturii comparativ cu cea obținută prin programul TimberWolfSC, și o reducere a timpului de execuție.

O posibilitate de rezolvare a problemei de plasare este transformarea acesteia într-o problemă de optimizare numerică. A fost descrisă o metodă numită plasare controlată de forțe, în care problema de plasare este redusă la soluționarea unui sistem de ecuații liniare pentru a determina locațiile de echilibru (coordonatele ideale x , y) ale celulelor. Dezavantajul acestor metodelor numerice este că ele necesită un timp de calcul foarte mare.

Metodele spectrale au fost utilizate și pentru problema de plasare. În cazul unei asemenea metode, prin utilizarea unei funcții liniare se obține o plasare de calitate mai bună din punct de vedere a lungimii conexiunilor. Utilizarea unei funcții obiectiv cuadratică are însă ca efect obținerea unui număr mai redus de conexiuni foarte lungi față de cazul unei funcții obiectiv liniare. A fost descrisă o metodă de plasare liniară propusă de Li care utilizează o funcție obiectiv spectrală constând dintr-un compromis între funcția cuadratică și cea liniară, ceea ce permite folosirea avantajelor oferite de ambele funcții.

Dintre metodele neconvenționale de plasare, a fost prezentată plasarea prin algoritmi paraleli și plasarea prin rețele neuronale artificiale. A fost descrisă implementarea paralelă a unui algoritm de plasare care se bazează pe metoda de călire simulată, care are aceleași proprietăți de convergență ca și algoritmul serial. Algoritmul utilizează două moduri de paralelizare, în funcție de valoarea temperaturii. A fost prezentată o soluție a problemei de plasare elaborată de Yu, care a utilizat rețelele neuronale Hopfield. Acesta a modificat soluția propusă de Hopfield și Tank pentru rezolvarea problemei comis-voiajorului.

În acest capitol s-au propus doi algoritmi de plasare pentru circuitele FPGA cu resurse limitate de rutare, în particular pentru circuitele FPGA *Atmel*. Acești algoritmi au ca obiectiv asigurarea rutabilității circuitelor. Primul algoritm se bazează pe algoritmul de partiționare pe baza tăieturii minime propus în capitolul 3. Algoritmul de plasare propus utilizează bipartiționarea a cărei funcție obiectiv urmărește minimizarea lungimii interconexiunilor simultan cu distribuția echilibrată a conexiunilor din cele două porțiuni.

A fost propusă de asemenea o secvență de aplicare a liniilor de tăietură care este mai eficientă decât secvențele tradiționale. În cadrul plasării convenționale bazate pe tăietura minimă, se alege linia de tăietură poziționată în centrul regiunii curente. Această alegere are ca rezultat o dimensiune a tăieturii relativ ridicată pentru liniile de tăietură din apropierea centrului, ceea ce are ca efect creșterea necesarului de resurse de rutare din această zonă. În secvența propusă, în zona centrală liniile de tăietură sunt aplicate în primele etape ale procesului de bipartiționare. Rezultatele experimentale arată că prin aplicarea acestei secvențe de tăietură se obține atât o reducere a dimensiunii maxime a tăieturii, cât și o reducere a sumei totale a tăieturilor, comparativ cu procedura de plasare quadratică.

Al doilea algoritm propus este un algoritm genetic pentru plasarea circuitelor FPGA, după cunoștințele noastre nefiind publicat până în prezent un algoritm genetic pentru plasarea acestor circuite. În cadrul algoritmului, strategia de selecție a părinților încearcă să evite convergența prematură a algoritmului la un optim local. În algoritmul propus, strategia de înlocuire este diferită de cea a algoritmilor tradiționali. După generarea urmașilor, aceștia înlocuiesc un număr echivalent de indivizi din populația curentă, aleși dintre cei cu viabilitatea cea mai redusă. Această strategie permite supraviețuirea soluțiilor mai bune de-a lungul unui mare număr de generații.

Pentru facilitarea rutării, algoritmul alocă un număr de celule libere în cadrul procesului de plasare, în scopul utilizării acestor celule pentru rutare. Algoritmul utilizează o funcție de cost care optimizează diferite metrici, care cuprind atât lungimea interconexiunilor, cât și măsuri ale rutabilității plasării. Funcția de cost cuprinde o componentă de estimare a rutabilității locale, cât și o componentă de echilibrare a densității. Algoritmul genetic a fost comparat cu un algoritm de plasare pe baza metodei de călire simulată, rezultatele experimentale arătând că în cazul algoritmului genetic timpul de execuție este mai redus, la o calitate a soluției care este apropiată de cea obținută prin metoda de călire simulată.

Contribuțiile acestui capitol sunt următoarele:

- Elaborarea și testarea unui algoritm de plasare pe baza bipartiționării, a cărei funcție obiectiv urmărește minimizarea lungimii interconexiunilor simultan cu distribuția echilibrată a conexiunilor din cele două porțiuni.
- Propunerea unei secvențe de aplicare a liniilor de tăietură care este mai eficientă decât secvențele tradiționale, având ca efect o reducere a dimensiunii maxime a tăieturii, cât și a sumei totale a tăieturilor, deci implicit a lungimii totale a interconexiunilor.
- Elaborarea și testarea unui algoritm genetic pentru plasarea circuitelor FPGA, având ca obiectiv atât reducerea lungimii totale a interconexiunilor, cât și asigurarea rutabilității circuitului.