

NIVERSITATEA "POLITEHNICA"
TIMIȘOARA

BIBLIOTECA CENTRALĂ

Nr. inv. 632.449

Dulap 145 Lit. C

SITATEA "POLITEHNICA" TIMIȘOARA
ltatea de Electronică și Telecomunicații

Ing. Lucian Aurel JURCA

PROCESOR ARITMETIC LOGARITMIC CU VIRGULĂ FLOTANTĂ

TEZĂ DE DOCTORAT

Conducător științific:

Prof. dr. ing. Mircea CIUGUDEAN

2002

Ing. Lucian Aurel JURCA

**PROCESOR ARITMETIC LOGARITMIC
CU VIRGULĂ FLOTANTĂ**

TEZĂ DE DOCTORAT

BIBLIOTECA CENTRALĂ
UNIVERSITATEA "POLITEHNICA"
TIMIȘOARA

Conducător științific:

Prof. dr. ing. Mircea CIUGUDEAN

2002

CUPRINS

Introducere.....	1
Cap.1 Metode de operare numerică.....	6
1.1 Circuite rapide de operare cu numere întregi.....	6
1.2 Metode de calcul utilizate în aritmetica binară în virgulă flotantă.....	14
1.2.1 Formate de date în virgulă flotantă.....	14
1.2.2 Realizarea celor patru operații aritmetice elementare în virgulă flotantă.....	15
1.2.3 Metode de creștere a vitezei operațiilor de adunare și scădere în virgulă flotantă.....	21
1.3 Metode de operare numerică în aritmetică logaritmică.....	24
1.3.1 Realizarea operațiilor de înmulțire și împărțire a două numere cu ajutorul logaritmilor.....	24
1.3.2 Realizarea unor operații matematice complexe cu ajutorul logaritmilor binari.....	27
1.3.3 Realizarea operațiilor de adunare și scădere în aritmetica logaritmică.....	29
Cap.2 Concepția subunității logaritmice a unui procesor aritmetic hibrid, ce permite efectuarea operațiilor de înmulțire și împărțire.....	35
2.1 Premize teoretice. Deziderate propuse.....	36
2.2 Subunitate de calcul logaritmică pentru realizarea rapidă a operațiilor de înmulțire și împărțire în simplă precizie.....	37
2.2.1 Algoritm de realizare a conversiilor de format FP-LNS și LNS-FP.....	37
2.2.2 Proiectarea subunității logaritmice după metoda Lai folosind algoritmul de conversie de format modificat.....	44
2.2.2.1 Determinarea dimensiunii memoriilor. Proiectarea circuitului de logaritmare.....	44
2.2.2.2 Proiectarea multiplicatorului de 12×12 biți.....	48
2.2.2.3 Circuitul de antilogaritmare.....	50
2.2.2.4 Proiectarea circuitului de adunare/scădere a operanzilor în format LNS cu deplasament.....	51
2.2.2.5 Concluzii la paragraful 2.2.2.....	58
2.2.3 Implementarea algoritmului modificat de conversie de format FP-LNS și LNS-FP prin metoda logaritmului mascat.....	60
2.2.3.1 Metoda logaritmului mascat.....	60
2.2.3.2 Prezentarea noului tip de compresor 4:2 folosit în proiectarea subunității logaritmice de calcul.....	67
2.2.3.3 Concepția blocului ALU în cazul aplicării metodei logaritmului mascat.....	72
2.2.3.4 Concluzii la paragraful 2.2.3.....	74

2.2.4 Studiul erorilor produse la generarea logaritmului și antilogaritmului binar.....	75
Cap.3 Concepția unei subunități de calcul în virgulă flotantă pentru realizarea operațiilor de adunare și scădere	78
3.1 Proiectarea circuitului de aliniere a mantiselor.....	79
3.2 Proiectarea circuitului de adunare/scădere a mantiselor aliniate.....	85
3.3 Proiectarea circuitului de normalizare	90
Cap.4 Rezultatele simulărilor blocurilor și subblocurilor componente ale procesorului hibrid în virgulă flotantă și logaritmic.....	94
4.1 Rezultatele simulărilor blocurilor și subblocurilor componente ale subunității de calcul logaritmice.....	95
4.2 Rezultatele simulărilor blocurilor și subblocurilor componente ale subunității de calcul în virgulă flotantă.....	104
Concluzii, contribuții și perspective.....	113
Bibliografie.....	117
Anexe	

Listă de abrevieri

- ESC – Prefix ce indică o instrucțiune specifică procesorului aritmetic (*Escape*);
FP – Sistem de reprezentare a numerelor în virgulă flotantă (*Floating-Point*);
FPU – Unitate de calcul în virgulă flotantă (*Floating-Point Unit*);
LNS – Sistem logaritmic de reprezentare a numerelor (*Logarithm Number System*);
NDP – Procesor de date numerice, sau coprocesor (*Numeric Data Processor*);
NPX – Extensie numerică a procesorului, sau coprocesor (*Numeric Processor Extension*);
CU – Unitatea de control a coprocesorului (*Control Unit*);
NEU – Unitatea de execuție numerică a coprocesorului (*Numeric Execution Unit*);
CLA – Sumator cu transport anticipat (*Carry Look-Ahead Adder*);
CSA – Sumator cu conservarea transportului (*Carry Save Adder*);
BS – Circuit de deplasare programabil (*barrel shifter*);
ALU – Unitate aritmetică și logică;
NZS – Numărător de zerouri semnifivative.

Introducere

Procesoarele din familia Intel 80×86, până la apariția procesorului 80486DX, sunt unități de prelucrare a numerelor întregi, după diferite formate. Oricare însă dintre CPU - uri care operează cu numere întregi poate efectua prin soft operații cu numere reale, emulatoarele de coprocesor fiind un bun exemplu în acest sens. Totuși utilizarea unui procesor dedicat operațiilor matematice, adică a unui coprocesor (ca de exemplu 8087, 80287 sau 80387, apărute între 1980 și 1986), a permis reducerea semnificativă a timpului de calcul.

Majoritatea programelor care utilizează coprocesorul și care deci beneficiază de precizia și viteza acestuia sunt aplicații profesionale și științifice foarte răspândite cum ar fi: programe de prelucrare a datelor științifice, aplicații grafice, aplicații financiare sau aplicații funcționând în timp real precum controlul mașinilor sau al proceselor, navigație, robotică, achiziții de date etc.

Începând din 1989, odată cu apariția procesorului Intel 80486DX, procesorul principal și coprocesorul au putut fi integrate pe același cip. Cu toate acestea, coprocesorul este totuși considerat ca și unitate de calcul distinctă. Alți termeni utilizați pentru acesta sunt conform [31]: procesor aritmetic, extensie numerică a procesorului (NPX), procesor de date numerice (NDP) sau procesor de calcul numeric, unitate în virgulă flotantă (FPU).

Acest ultim termen, stabilit prin extensie, nu este întotdeauna corect. Tipul unui coprocesor este dat de modul în care sunt reprezentate datele în interiorul său. Atunci când sunt cerute simultan o scară largă de reprezentare a numerelor reale cât și precizie, se adoptă în general standardul sau sistemul de reprezentare în virgulă flotantă (FP). Încă de la început unitățile de calcul în virgulă flotantă au oferit suficiente avantaje pentru a fi larg răspândite și de aceea performanțele lor au fost în mod continuu îmbunătățite (ca de exemplu coprocesoarele integrate produse de firma Intel). Creșterea densității de integrare însă a permis dezvoltarea ca și alternativă și a unităților de calcul logaritmice [2], [4], [12], [18], [21], [22], [37], [40], [42], [44], [45], [50], [51], [57], [58], [59], [89], [90], [104]. Proprietățile aritmeticii logaritmice permit dezvoltarea unor algoritmi de calcul competitivi, care în cazuri particulare oferă o viteză de operare mai mare decât sistemele de calcul în virgulă flotantă. În același timp, precizia pe care o oferă sistemul LNS (Logarithm Number System) în reprezentarea și manipularea numerelor, este cu ceva mai mare decât a sistemului FP [57], pentru aceeași dimensiune a formatului. Având în vedere cele două aspecte, în lucrarea de față termenii FPU sau procesor LNS nu vor fi folosiți pentru a desemna la modul general un coprocesor sau procesor aritmetic. În schimb va fi prezentată structura unui coprocesor hibrid care va avea calități superioare oricăruia dintre cele două tipuri menționate.

Acest coprocesor va implementa cele patru operații aritmetice de bază și va fi conceput astfel încât să fie foarte eficient în efectuarea unor sume/diferențe multiple de produse/fracții. Așadar el va putea fi folosit în principal ca și procesor de calcul,

specializat, putând implementa cu ușurință prin microprogramare diverși algoritmi DSP (procesor de semnal), sau algoritmi pentru transformări de coordonate în cazul unor aplicații grafice. Având în vedere acest aspect, unitatea de execuție numerică ce va fi prezentată în această lucrare va putea fi integrată într-un cip de sine stătător, în urma asocierii unui hardware care să controleze un set specializat de instrucțiuni.

În același timp, vor fi indicate soluțiile care să permită extinderea cu ușurință a capabilităților sale de operare rapidă și în cazul calculului unor funcții mai complexe cum ar fi rădăcina pătrată, ridicarea la putere, exponențiala și logaritmul în orice bază. În acest context, coprocesorul hibrid va putea fi integrat și interfațat, în cadrul aceluiași cip, cu procesorul principal, putând suporta un set lărgit de instrucțiuni prin care să se efectueze operații aritmetice simple sau calculul unor funcții matematice complexe (cum ar fi cele trigonometrice sau transcendente) dintre cele mai diverse.

Setul de instrucțiuni al coprocesorului va utiliza o codare similară cu a procesorului. Instrucțiunile care controlează coprocesorul sunt o extensie simplă a limbajului obișnuit al procesorului și se numesc instrucțiuni "escape" (ESC) [31]. Cum procesorul și coprocesorul sunt două unități de calcul distincte (chiar dacă sunt integrate pe același cip), coprocesorul este capabil să lucreze, într-un mod asincron, în paralel cu procesorul. Când procesorul întâlnește o instrucțiune ESC (deci destinată coprocesorului), el calculează adresa memoriei dacă un operand este specificat și trece imediat să decodeze și să execute instrucțiunea următoare. Dacă aceasta nu este tot de tip ESC și nici nu utilizează rezultatul instrucțiunii ESC în curs de execuție, cele două procesoare lucrează în paralel. O codare îngrijită în asamblor (ținând cont că o operație mai complexă poate dura sute de clock-uri) permite folosirea paralelismului pentru a obține performanțe îmbunătățite.

Operanzii sunt situați fie în memorie, fie în regiștrii coprocesorului, iar adresarea memoriei se face în același mod ca în cazul procesorului. Setul de instrucțiuni cuprinde instrucțiuni suplimentare de încărcare și salvare pentru a transmite operanzi reali, întregi și BCD între regiștrii coprocesorului și memoria sistemului.

Indiferent de tipul sau destinația unui procesor aritmetic, în cadrul lui pot fi identificate două unități distincte: unitatea de control (CU) și unitatea de execuție numerică (NEU).

Unitatea de control efectuează interfațarea cu magistrala și decodarea instrucțiunilor coprocesorului, accede și reține operanzii și opcodurile (cod mașină, cod executabil sau cod obiect) și le transmite apoi către NEU. CU este responsabilă de sincronizarea coprocesorului cu procesorul. Sincronizarea este necesară când cele două procesoare trebuie să schimbe informații, să accedă la rezultatele operațiilor sau să decodeze instrucțiuni.

Unitatea de control este legată la liniile de stare ale procesorului și este capabilă să aducă instrucțiuni în paralel cu procesorul, atunci când acesta le citește din memorie. Coprocesorul trebuie să citească și să decodeze aceleași instrucțiuni

Introducere

ca și procesorul principal; trebuie deci să emuleze precitirile și decodarea instrucțiunilor proprii procesorului asociat. De exemplu, coprocesorul trebuie să mențină o coadă de așteptare de precitare de aceeași mărime cu a procesorului, să o manipuleze în așa fel încât instrucțiunile să fie plasate sau retrase (cu ocazia inițializărilor) în același timp cu procesorul.

Instrucțiunile destinate coprocesorului sunt inserate direct în codul destinat procesorului. Astfel procesorul și coprocesorul decodează împreună clasa de instrucțiuni ESC. Procesorul decodează suficient instrucțiunile ESC pentru a diferenția pe cele care fac referire la memorie de cele care nu fac. Dacă se face referire la memorie, procesorul calculează adresa operandului și începe un ciclu de citire pe magistrală, dar ignoră data primită. Dacă nu se face referire la memorie, procesorul începe imediat să decodeze și să execute instrucțiunea următoare. Procesorul ignoră operațiile incluse în instrucțiunile ESC.

În caz de acces la memorie, NPX salvează adresa efectivă a operandului. Dacă un operand registru este specificat, valoarea acestui operand este salvată. Dacă extragerea unui operand durează mai mult de un ciclu de magistrală (un transfer de date) din cauza mărimii sau aliniamentului său, CU a coprocesorului ia controlul magistralei sistemului și transferă restul operandului. La scrierea unui operand în memorie, după calculul adresei de către procesor, NPX salvează adresa și începe ciclul de scriere în memorie a operandului.

Există instrucțiuni ale coprocesorului executate numai de CU (fără participarea NEU) care indică sau modifică starea coprocesorului, ca de exemplu: FINIT (inițializarea NPX), FCLEX (ștergerea fanioanelor de excepție), FSTSWAX (stocarea cuvântului de stare în registrul AX al CPU), FSTCW (stocarea cuvântului de control într-un operand în memorie).

Unitatea de execuție numerică (NEU) este veritabilul coprocesor.

Aceasta efectuează toate operațiile prin care se accede la regiștrii numerici ai NPX sau prin care se prelucrează datele conținute de aceștia. Efectuează operații aritmetice, logice, transcendente, furnizând și un anumit număr de constante matematice înscrise într-un ROM intern al cipului. NEU conține circuite specializate precum sumatoare rapide, circuite combinate de adunare/scădere, multiplicatoare rapide, circuite pentru împărțire și de extragere a rădăcinii pătrate, circuite de deplasare programabile, circuite combinaționale dedicate (cum ar fi circuitele de contorizare a numărului de zerouri de după virgulă, sau de anticipare a numărului acestora etc.) utilizate după nevoi pentru efectuarea unor funcții matematice specifice. Pe măsură ce instrucțiunile ESC sunt decodate, ele trec în secvențiatorul de instrucțiuni care programează subunitățile cipului pentru realizarea funcției cerute.

Când execuția instrucțiunii ESC începe, NEU activează semnul său BUSY (ocupat) care poate fi supravegheat de procesorul principal, sau această stare este raportată în cuvântul de control al coprocesorului. Semnalul BUSY, combinat cu

instrucțiunea WAIT a procesorului principal permite resincronizarea procesorului și a coprocesorului său atunci când este necesar.

În lucrarea de față va fi oferită o soluție proprie privind proiectarea unității de execuție numerică (NEU) care să permită concepția unui procesor aritmetic performant. Acesta, fie va primi o destinație generală și va fi atașat unui procesor principal (caz în care vor fi indicate posibilitățile de extensie a capacităților NEU), fie va reprezenta o unitate de calcul specializată, de sine stătătoare, dedicată unui anumit context matematic, respectiv dedicată anumitor tipuri de aplicații.

Ca punct de pornire, în capitolul 1 al tezei vor fi analizate principalele metode de operare numerică și în primul rând operarea cu numere întregi, întrucât oricare altă metodă se reduce în cele din urmă la operare (adunare) de numere întregi. Vor fi astfel prezentate și analizate principalele tipuri de sumatoare rapide cunoscute în literatura de specialitate, realizate într-o astfel de manieră încât informația numerică ce se prelucrează să străbată un număr cât mai mic de niveluri logice. În următoarele două paragrafe vor fi prezentați algoritmi prin care se implementează diversele operații aritmetice, atunci când informația binară este reprezentată în virgulă flotantă (FP) în format simplă sau dublă precizie, respectiv în virgulă fixă, în format logaritm (LNS) simplă precizie. Vor fi evidențiate avantajele și dezavantajele fiecărui mod de reprezentare și de operare a informației numerice, ca și premise teoretice pentru concepția unui nou tip de coprocesor.

Pe baza concluziilor desprinse în primul capitol, în capitolul 2, cel mai important al tezei, va fi descrisă structura unui coprocesor hibrid, care va îngloba o subunitate logaritmă de calcul, ce va avea intrări și ieșiri în formatul virgulă flotantă, simplă precizie. Capitolul 2 este aproape în întregime original, motiv pentru care el este redactat în mare parte cu caractere *italice* pentru a evidenția contribuțiile autorului. De asemenea, figurile și relațiile originale, precum și citarea lucrărilor proprii vor fi notate cu caractere *italice*. Subunitatea logaritmă va include convertoare de format la intrarea și ieșirea sa, care să permită operarea în interiorul său în format LNS, în scopul realizării rapide a operațiilor de înmulțire și împărțire. De asemenea va fi prezentat un studiu al erorilor rezultate după efectuarea conversiilor de format și va fi prezentată o metodă de reducere a acestora, în condițiile în care și aria ocupată pe cip va fi mai mică.

În capitolul 3 va fi concepută o subunitate în virgulă flotantă, pentru implementarea operațiilor de adunare și scădere. Structura proiectată va fi optimizată din punct de vedere al suprafeței ocupate pe cip și al realizării sincronizării cu subunitatea logaritmă. Astfel, ansamblul format din cele două subunități va funcționa în parametri de mare performanță, mai ales în cazul implementării prin microprogramare a unor algoritmi DSP, sau transformări de coordonate pentru aplicații grafice.

În capitolul 4 va fi demonstrată viteza mare de calcul a tuturor structurilor proiectate, pe baza simulărilor efectuate cu ajutorul programului "MSim" în cazurile

Introducere

cele mai defavorabile din punct de vedere al propagării transportului prin toate aceste circuite.

În continuare, în ultimul capitol al tezei, vor fi prezentate concluziile desprinse în urma rezultatelor obținute, vor fi evidențiate toate contribuțiile teoretice și aplicative aduse de autor pe cuprinsul lucrării și vor fi enunțate perspectivele de dezvoltare ulterioară a temei abordate.

În anexele de la lucrare va fi prezentată pe părți componente schema completă a unității de execuție numerică a coprocesorului hibrid, programele Matlab pentru generarea valorilor ce trebuie memorate în ROM, în condițiile minimizării erorii de conversie de format, programul C pentru depistarea cazurilor cele mai defavorabile pentru sumatorul cu transport anticipat pe 8 biți, precum și un exemplu de rulare al unui algoritm DSP simplu, prin care se efectuează o buclă Livermore cu 20 de termeni (suma a 20 de produse de câte doi factori).

Cap.1 Metode de operare numerică

1.1 Circuite rapide de operare cu numere întregi

Indiferent de tipurile de date cu care lucrează coprocesorul, în final orice operație, indiferent cât de complexă ar fi ea, se reduce la o operație cu numere întregi. Aceasta înseamnă că pentru a îmbunătăți viteza de calcul într-un procesor aritmetic trebuie acționat inclusiv în această direcție, respectiv trebuie alese acele structuri de operare cu numere întregi care oferă cei mai mici timpi de producere a rezultatului, iar pentru structuri care furnizează timpi comparabili, se alege acel circuit care ocupă cea mai mică suprafață pe cipul de siliciu. În funcție de lungimea cuvintelor de date cu care se operează, se poate depista prin analiză și simulare care este tipul optim de structură ce trebuie adoptat, pentru o tehnologie dată.

Cea mai frecvent utilizată operație elementară este operația de adunare. Toate celelalte operații elementare se reduc sau cel puțin fac apel și la operația de adunare. Astfel, operația de scădere se efectuează prin adunarea la descăzut a complementului de doi al scăzătorului ($A-B=A+/\overline{B}+1$, în care $/\overline{B}$ este \overline{B} negat). Operația de înmulțire se realizează prin însumarea produselor parțiale iar operația de împărțire include scăderi repetate întreșesute cu comparări.

De asemenea, calculul funcțiilor trigonometrice și transcendente se face prin intermediul unor serii rapide sau algoritmi, [22], [23], [48], [88], al căror calcul implică tot utilizarea celor patru operații aritmetice elementare. Aceasta înseamnă că practic sumatoarele binare de numere întregi sunt printre cele mai solicitate componente hard într-un coprocesor, motiv pentru care ele vor fi prezentate și analizate în acest paragraf.

Sumatorul simplu

Cel mai simplu sumator binar a două numere de n biți, dar și cel mai dezavantajos din punct de vedere al timpului de propagare, este cel realizat prin cascadarea a n sumatoare complete pe un bit. Un astfel de sumator complet (sumator complet 3:2) generează pe un bit suma s_i a celor 2 biți de rang i din cei doi operanzi ținând cont de transportul c_i furnizat de sumatorul anterior, iar pe un alt bit furnizează transportul c_{i+1} către sumatorul următor. Pot fi scrise următoarele relații:

$$s_i = a_i \oplus b_i \oplus c_i \qquad c_{i+1} = a_i \cdot b_i + b_i \cdot c_i + c_i \cdot a_i \qquad (1.1)$$

Aceste relații conduc la implementarea prezentată în figura 1.1. Întrucât $a \oplus b = a \cdot \overline{b} + \overline{a} \cdot b$ se obține pentru s_i următoarea expresie:

$$s_i = a_i \cdot \overline{b_i} \cdot \overline{c_i} + \overline{a_i} \cdot b_i \cdot \overline{c_i} + a_i \cdot b_i \cdot c_i + \overline{a_i} \cdot \overline{b_i} \cdot c_i \qquad (1.2)$$

Deoarece transportul se propagă de la o celulă la alta traversând câte două niveluri logice (conform relației 1.2), rezultă că e nevoie de $2n$ niveluri logice pentru obținerea rezultatului. Pentru creșterea vitezei operației de adunare au fost adoptate

diverse tehnici care să permită reducerea numărului de niveluri logice parcurse de informația binară ce se prelucrează.

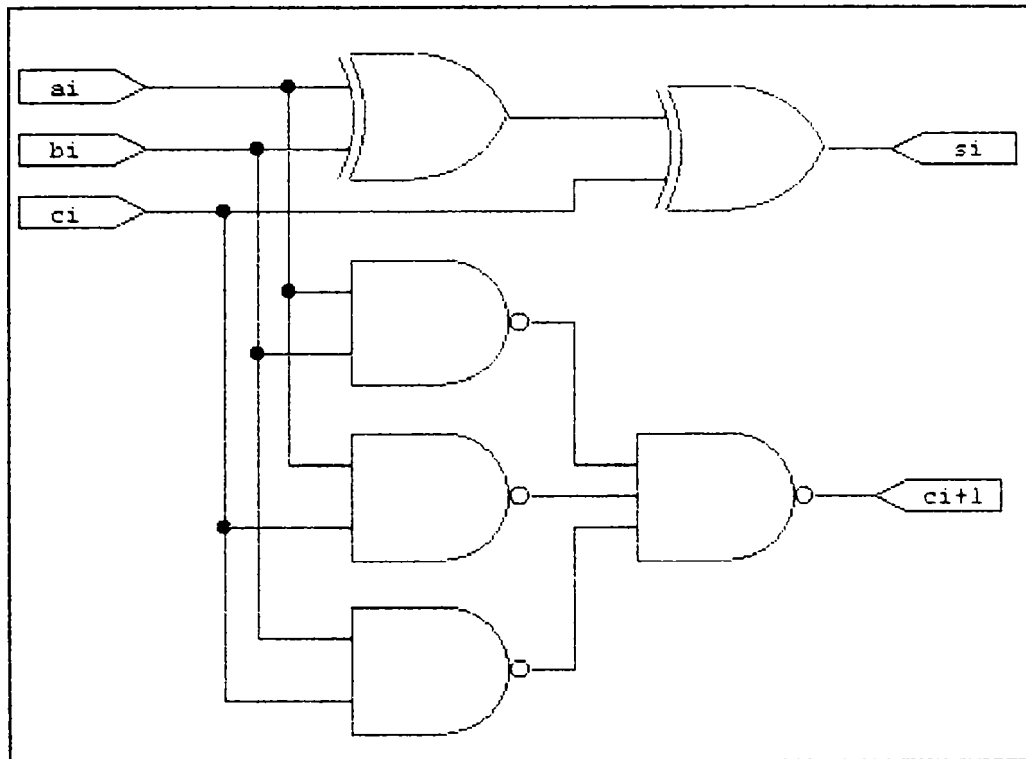


Fig. 1.1 Schema unui sumator complet pe un bit (sumator complet 3:2)

Reducerea numărului de niveluri se poate face utilizând următoarele relații:

$$c_{i+1} = g_i + p_i c_i \quad \text{în care:} \quad g_i = a_i b_i \quad \text{și} \quad p_i = a_i + b_i \quad (1.3)$$

Făcând înlocuirile din aproape în aproape obținem:

$$c_{i+1} = g_i + p_i g_{i-1} + p_i p_{i-1} g_{i-2} + \dots + p_i p_{i-1} \dots p_1 g_0 + p_i p_{i-1} \dots p_1 p_0 c_0 \quad (1.4)$$

În concluzie e nevoie de un nivel logic pentru a forma termenii p și g , două niveluri logice pentru a forma transportul cu relația 1.4 și două niveluri logice pentru sumă. În total rezultă cinci niveluri logice. Dezavantajele structurii care ar implementa relațiile 1.3 și 1.4 constau în faptul că sunt necesare porți cu foarte multe intrări (de exemplu, semnalul p_{n-1} trebuie să apară la intrarea a n porți ȘI), rezultând astfel o structură neregulată și foarte multe fire de legătură. Această idee totuși a fost folosită pentru a se construi un sumator care va avea $\log_2 n$ niveluri logice, deci mult mai puține decât un sumator cu propagarea transportului și în plus o structură mai simplă și mai regulată.

Sumatorul cu transport anticipat

Sumatorul cu transport anticipat este cunoscut în literatura de specialitate, [73], [80], sub denumirea de “carry look-ahead adder” (CLA). Ideea folosită constă în construirea semnalelor auxiliare g și p în pași. Există deja $c_1 = g_0 + c_0 p_0$, care înseamnă că un transport este generat din poziția 0 fie datorită faptului că poziția în sine (prin a_0 și b_0 egali cu 1) generează transport, fie datorită faptului că un transport anterior se poate propaga peste această poziție (dacă a_0 sau b_0 este egal cu 1).

Similar $c_2 = G_{01} + P_{01}c_0$ unde $G_{01} = g_1 + p_1g_0$ și $P_{01} = p_1p_0$. G_{01} arată că blocul format din primii doi biți generează un transport, iar P_{01} arată că prin prin acest bloc se poate propaga un transport anterior.

În general, considerând $i < j < k$ putem scrie următoarele relații:

$$C_{k+1} = G_{ik} + P_{ik}c_i \tag{1.5}$$

$$G_{ik} = G_{j+1,k} + P_{j+1,k}G_{ij} \tag{1.6}$$

$$P_{ik} = P_{ij} P_{j+1,k} \tag{1.7}$$

în care $P_{ii} = p_i$ și $G_{ii} = g_i$. Relația (1.5) se interpretează astfel: un transport este generat la ieșirea blocului format din biții de la i la k inclusiv, dacă partea superioară a acestui bloc (incluzând biții de rang cuprins între $j+1$ și k) generează un transport, sau dacă există un transport generat de partea inferioară (incluzând biții de rang cuprins între i și j) care se poate propaga prin partea superioară.

Pe baza acestor considerații se prezintă o variantă practică de sumator CLA [80]. Sumatorul constă din două părți. Prima parte (figura 1.2) calculează diverse valori P și G din p_i și g_i utilizând ecuațiile (1.6) și (1.7).

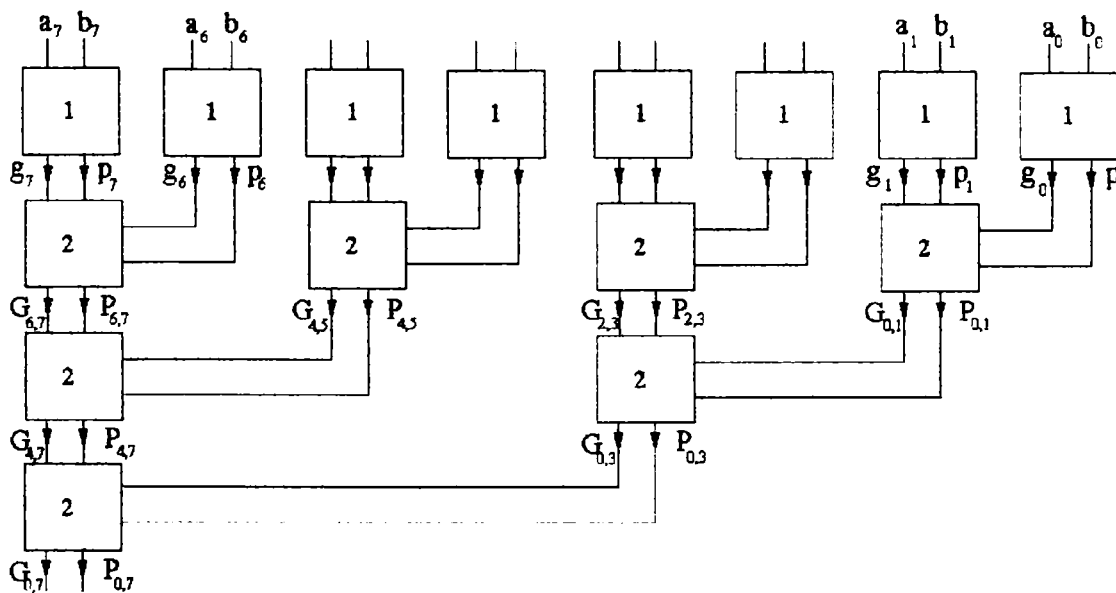


Fig.1.2 Schema bloc a circuitului care generează coeficienții G și P

Celulele de tip 1 vor calcula valorile g_i și p_i din a_i și b_i , iar celulele de tip 2 interconectate într-o structură de tip arbore, grupate două câte două, vor calcula perechile corespunzătoare (P, G).

A doua parte folosește valorile P și G pentru a calcula toate transporturile cu ajutorul ecuației (1.5). Schema care se obține este prezentată în figura 1.3.

Cele două tipuri de celule pot fi alăturate într-o structură unică așa cum se vede în figura 1.4. Astfel, într-un sumator cu structură arborescentă de tip CLA, numerele binare ce trebuie adunate intră în partea de sus a arborelui, curg în jos pentru a se combina cu c_0 și apoi informația curge în sus pentru a se calcula biții sumei. Numărul de celule utilizate este $2n$ (inclusiv cea necesară pentru generarea lui c_8), deci de două ori mai multe decât la un sumator clasic, însă din punct de

vedere al spațiului ocupat pe cip suprafața ocupată este egală cu $n \times \log_2 n$ celule dintr-un sumator clasic.

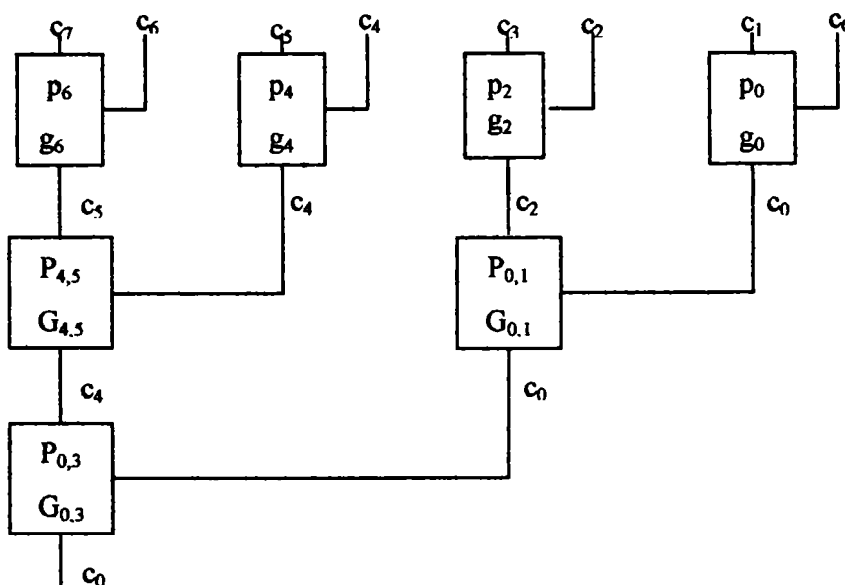


Fig. 1.3 Schema bloc a circuitului care generează transporturile c_i

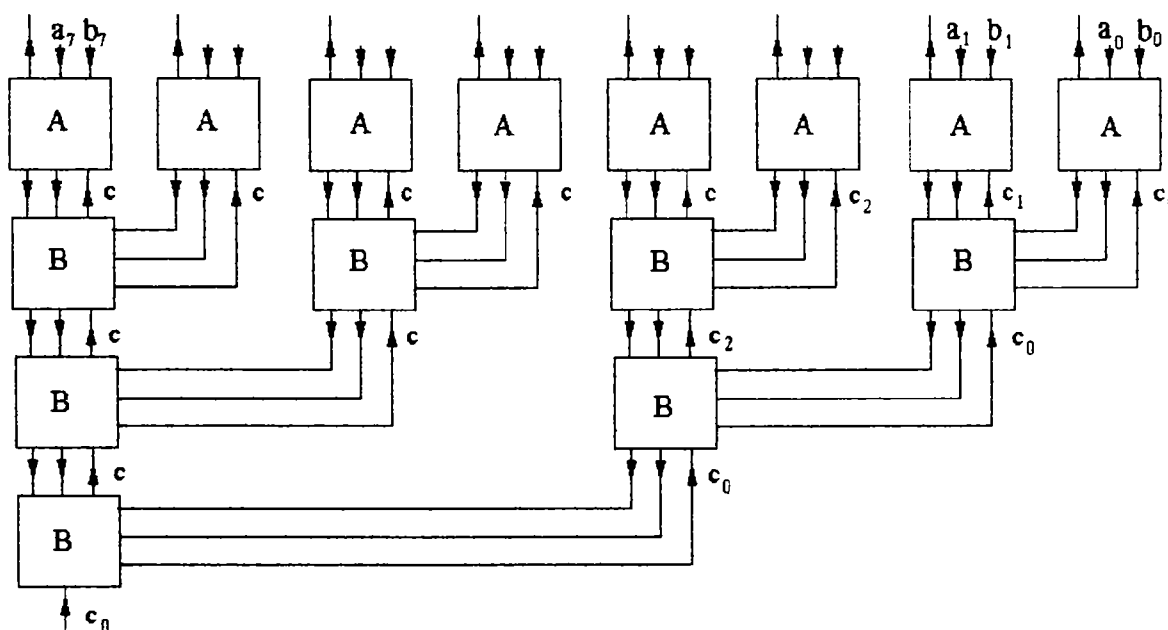


Fig.1.4 Schema bloc a unui sumator CLA pe 8 biți, obținută prin suprapunerea structurilor prezentate în figurile 1.2 și 1.3

Structura blocurilor de tip A și B este prezentată în figurile 1.5.a și 1.5.b.

Concluzia care se desprinde totuși aici este că printr-o mică investiție în mărimea structurii se obține o creștere substanțială a vitezei, comparativ cu un sumator cu propagarea din aproape în aproape a transportului. Acest lucru este cu atât mai pregnant cu cât numărul n de biți ai numerelor ce trebuie adunate este mai mare.

Există mai multe tipuri de sumatoare CLA care se adoptă în funcție de tehnologia utilizată. De exemplu dacă fiecare nod al arborelui ar avea trei intrări în loc de două, atunci înălțimea arborelui ar descrește de la $\log_2 n$ la $\log_3 n$, dar și

celulele ar fi mai complexe. Pentru tehnologii în care propagarea semnalelor este rapidă, o soluție hibridă între un sumator cu propagarea transportului și un arbore CLA, ar fi un bun compromis. În acest caz blocurile de tip A nu ar grupa perechi de biți ci perechi de blocuri de biți și vor fi constituite din sumatoare simple cu număr redus de biți, blocurile fiind organizate într-un arbore CLA, rezultând pe ansamblu o structură cu anticiparea transportului pe două niveluri ierarhice [13], [78].

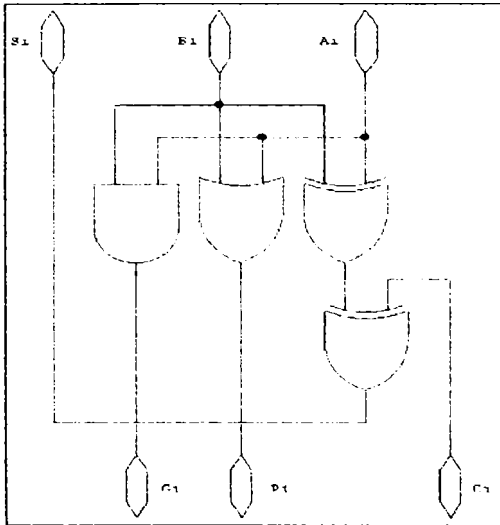


Fig. 1.5.a
Schema blocului de tip A

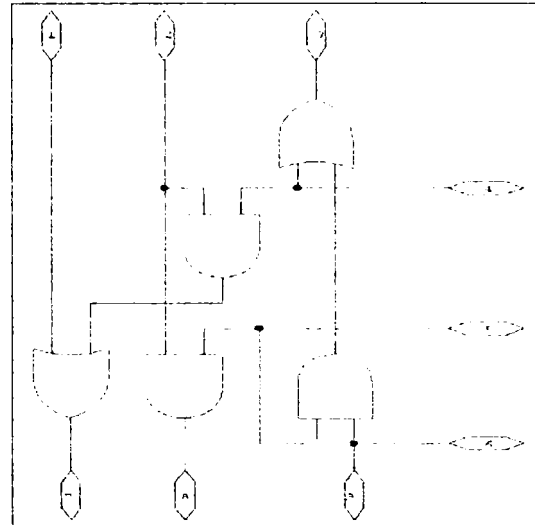


Fig. 1.5.b
Schema blocului de tip B

Sumatorul cu omiterea transportului

Un astfel de sumator se situează din punct de vedere al costului și al performanțelor între un sumator cu propagarea transportului și un sumator CLA. Un exemplu de sumator cu omiterea transportului este prezentat în figura 1.6. El este format din blocuri de sumatoare complete pe un bit cascade, notate B_i în figura 1.6. Numărul sumatoarelor complete din fiecare bloc se stabilește astfel încât să se obțină o viteză maximă de producere a rezultatului, pentru o lungime dată a operanzilor. Grupând de exemplu câte patru biții cuvintelor ce trebuie adunate și scriind pentru prima grupare ecuațiile coeficienților G și P se obține:

$$P_{03} = p_0p_1p_2p_3; \quad G_{03} = g_3 + p_3g_2 + p_3p_2g_1 + p_3p_2p_1g_0 \quad (1.8)$$

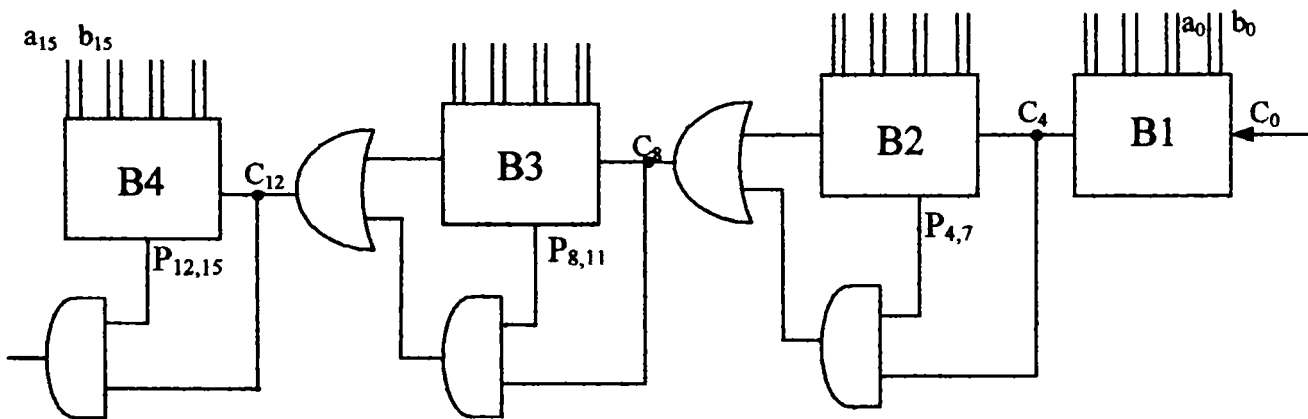


Fig. 1.6 Schema bloc a unui sumator cu omiterea transportului

Deoarece se observă că e mai simplu de calculat coeficienții P decât G, acest tip de sumator calculează doar coeficienții P. Propagarea transportului începe simultan în fiecare bloc. Dacă fiecare bloc va genera un transport atunci transportul generat de bloc va apare chiar dacă transportul de intrare al blocului nu este încă corect. Dacă la startul fiecărei operații de adunare, transportul la intrarea fiecărui bloc este zero, atunci nici un transport fals nu va fi generat, iar transportul generat poate fi apreciat ca fiind coeficientul G al acestuia.

După ce cel mai puțin semnificativ bloc își generează transportul, acesta va apare atât la intrarea blocului următor cât și la intrarea porții ȘI împreună cu semnalul P al aceluși bloc. Dacă atât semnalul P al acestui bloc cât și transportul de la blocul anterior sunt ambele 1 atunci imediat va fi generat transportul de intrare pentru următorul bloc. Astfel, a fost omis transportul pe care l-ar fi produs blocul actual în absența unui transport la intrarea sa.

Dimensiunea și numărul optim de blocuri este în funcție de lungimea cuvintelor de date care se adună, așa cum se arată în [47]. De asemenea fiecare bloc al sumatorului poate fi organizat sub forma unui sumator cu omiterea transportului, rezultând pe ansamblu un sumator cu omiterea transportului pe 2 niveluri ierarhice, ce prezintă o viteză de calcul substanțial mărită, [28], [46], comparabilă cu a unui sumator cu anticiparea transportului.

Sumatorul cu selecția transportului

Sumatorul cu selecția transportului funcționează după următorul principiu: două sumatoare simple lucrează în paralel în fiecare bloc al sumatorului (excepție făcând cel mai puțin semnificativ bloc), un sumator considerând transportul de intrare ca fiind "0", alta considerându-l 1. Când transportul de intrare este în final cunoscut, se selectează suma corectă.

În figura 1.7 este construit un sumator pe 19 biți format din patru blocuri. Dacă un bloc necesită k unități de timp pentru a calcula k biți și o unitate de timp pentru a furniza "carry in" al blocului următor din două semnale "carry out" disponibile și un "carry in" al blocului anterior, rezultă că pentru a optimiza operația de adunare e necesar ca fiecare bloc să fie mai larg cu 1 bit decât anteriorul.

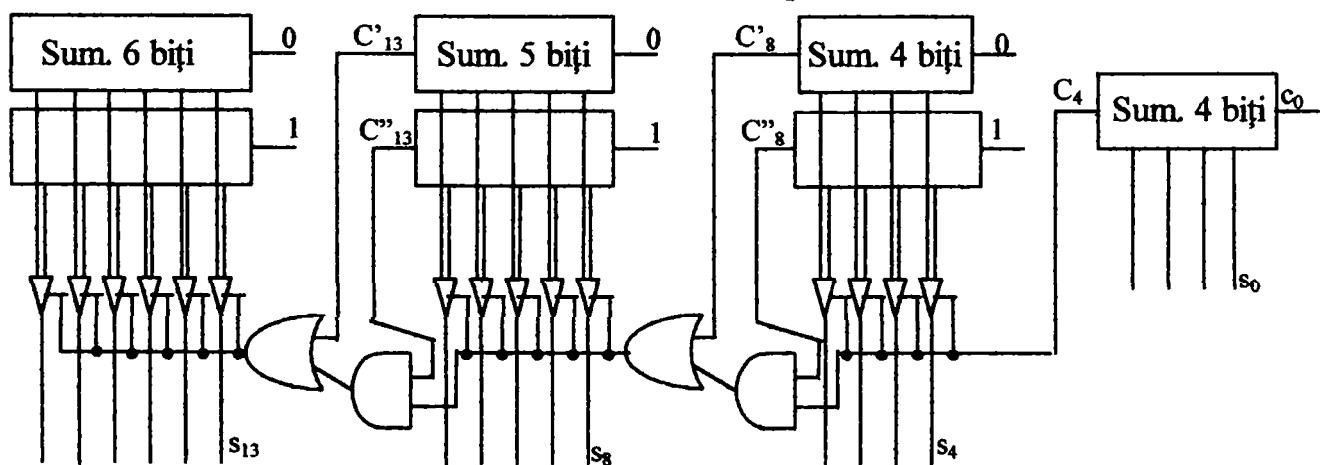


Fig. 1.7 Schema bloc a unui sumator cu selecția transportului

Și acest tip de sumator poate fi organizat pe două niveluri ierarhice în scopul creșterii vitezei de calcul. Aceasta înseamnă că în locul fiecărui bloc de sumatoare simple cascadeate poate fi introdus un sumator cu selecția transportului care va conține subblocuri cu sumatoare simple cascadeate.

În continuare va fi făcută o analiză a timpului de producere a rezultatului și a ariei ocupate pe cip, pentru tipurile de sumatoare prezentate anterior. Se consideră unitatea de timp de propagare, T , ca fiind timpul necesar pentru producerea transportului în cazul unui sumator complet pe 1 bit. Așa cum se observă în figura 1.1, c_{i+1} se obține după parcurgerea a două niveluri logice, deci după un timp T ; pentru obținerea bitului sumă s_i este necesar un timp $2T$. De asemenea se consideră unitatea de suprafață, S , ca fiind aria ocupată de un sumator complet pe 1 bit.

Astfel, pentru sumatorul simplu care adună cuvinte de n biți, timpul de calcul al rezultatului este $t=(n-1) \times T+2T=(n+1) \times T \approx n \times T$, dacă n este suficient de mare. Suprafața ocupată va fi $n \times S$.

Pentru sumatorul cu transport anticipat, ținând cont că transportul se propagă de sus în jos prin structura din figura 1.2 și apoi de jos în sus prin structura din figura 1.3 și că parcurgerea într-un sens sau altul a unui bloc de tip B (fig. 1.5.b) se face într-un timp T , iar a blocurilor de tip A (fig.1.5.a) în $0,5T$ (în jos), respectiv T (în sus), rezultă că timpul total de calcul este:

$$t_{CLA}=0,5 \times T+(\log_2 n-1) \times T+(\log_2 n) \times T+T=(2 \log_2 n) \times T+0,5 \times T \approx (2 \log_2 n) \times T \quad (1.9)$$

Am constatat astfel că formula prezentată în [80] ce furnizează t_{CLA} este greșită, ea furnizând o valoare ce reprezintă doar jumătate din valoarea corectă. Aria ocupată este $n \times \log_2 n \times S$ deoarece arborele CLA conține $\log_2 n$ etaje. Deoarece însă numărul de blocuri de tip B se înjumătățește spre vârful de jos al arborelui, o proiectare de layout optimă poate conduce la obținerea unei arii mai mici, fără a afecta caracterul modular al sumatorului.

În cazul sumatorului cu omiterea transportului, timpul total de calcul este dat de timpul de propagare al blocului cu cea mai mare dimensiune [47]. Acest tip de sumator este însă din ce în ce mai puțin utilizat deoarece el este eficient doar dacă se realizează în tehnologii în care semnalele pot fi șterse la începutul fiecărei operații [73] și care să permită integrarea de porți logice cu multe intrări [46].

Pentru sumatorul cu selecția transportului, numărul optim de blocuri componente, în cazul unei structuri regulate în care toate blocurile au aceeași dimensiune, este $n^{1/2}$ [73]. Dimensiunea unui astfel de bloc va rezulta $n^{1/2} \times S$, iar timpul de propagare al transportului prin bloc va fi deci $n^{1/2} \times T$. Considerând timpul necesar selecției rezultatului corect al unui bloc, de valoare cel puțin T , rezultă că timpul total de calcul pentru acest sumator este cel puțin $t_{SEL}=2n^{1/2} \times T$. Acest rezultat corectează valoarea greșită furnizată de [80]. Suprafața ocupată de acest sumator este dată de relația (1.10):

$$S_{SEL} = \left(\frac{\sqrt{n}-1}{\sqrt{n}} + 1 \right) n \times S + S_{\text{mux}} \approx 2n \times S, \quad (1.10)$$

în care S_{mux} este aria ocupată de circuitele multiplexoare care selectează rezultatul corect la nivelul fiecărui bloc.

Utilizarea unuia sau altuia dintre tipurile de sumatoare prezentate (ce reprezintă tipuri de bază, fiind printre cele mai des utilizate), depinde de locul și scopul în care sunt utilizate. Dacă ele apar în proiectare pe o cale critică de propagare (caz în care viteza lor influențează direct frecvența de clock a procesorului) atunci se folosesc cele mai rapide, chiar dacă aria pe care o vor ocupa pe cip este mare. Dacă nu este așa, atunci se alege tipul de sumator care ocupă aria cea mai mică, dar care este suficient de rapid pentru a nu deturna calea critică de propagare în structura proiectată.

Principalele caracteristici ale sumatoarelor prezentate, referitoare la timpul necesar pentru producerea rezultatului și la aria ocupată pe cip, sunt prezentate în tabelul 1.1.

Tab. 1.1

Unitate de timp de propag., T și de arie ocup., S	timp	supraf.
Sumator simplu cu n sum. complete cascade	$n \times T$	$n \times S$
Sumator cu anticiparea transportului	$2 \log_2 n \times T$	$n \times \log_2 n \times S$
Sumator cu omiterea transportului	$2n^{1/2} \times T$	$n \times S$
Sumator cu selecția transportului	$2n^{1/2} \times T$	$2n \times S$

Cu aceste tipuri de sumatoare se pot realiza diverse combinații care permit obținerea unor sumatoare mixte, mai rapide decât cele simple, organizate pe două niveluri ierarhice. Nivelul inferior se referă la structura blocurilor componente iar nivelul superior se referă la structura de ansamblu a sumatorului, respectiv la modul în care sunt interconectate blocurile componente.

Așa cum se va vedea pe parcursul lucrării, anumite combinații vor conduce fie la obținerea unei viteze superioare oricăruia din tipurile de bază menționate, fie la obținerea unui anumit optim, pe o cale necritică de propagare. Structura optimă pentru un anumit sumator (ce operează cu date de o lungime dată) se obține pentru acea configurație care să poată fi integrată pe cip într-o arie minimă și care să nu furnizeze rezultatul într-un timp mai mare decât cel impus. Motivația obținerii acestei structuri optime constă în aceea că nu are rost să se adopte un sumator foarte rapid, adică un circuit care ocupă o arie relativ mare pe cip, dacă el nu poate contribui la creșterea frecvenței de clock a sistemului de calcul.

În cazul în care unitatea de calcul este organizată pe niveluri pipeline, pe nivelul sau nivelurile critice de propagare (ce prezintă cel mai mare timp de producere a rezultatului și care impun practic frecvența de clock a întregului sistem) vor fi concepute sau adoptate cele mai rapide structuri numerice, economia de arie pe cip ne mai fiind în acest caz prioritară.

1.2 Metode de calcul utilizate în aritmetica binară în virgulă flotantă

1.2.1 Formate de date în virgulă flotantă

Într-un calculator numărul cifrelor disponibile pentru a reprezenta un număr real este limitat. Dacă avem de-a face cu un coprocesor, numărul de biți cu care sunt reprezentate numerele reale este însă suficient de mare pentru a avea o bună aproximație a acestora.

Atunci când sunt cerute simultan atât precizie cât și o scară largă de reprezentare a numerelor reale se apelează de obicei la modul de reprezentare în virgulă flotantă (FP). Formatele în virgulă flotantă conțin trei câmpuri: de semn, de exponent și de mantisă, așa cum se vede în figura 1.8. Numerele cu care se lucrează sunt, implicit, normalizate. Un număr binar este normalizat atunci când are un singur 1 la stânga virgulei. Excepția este numărul 0 care este tratat ca un caz special. Numărul 1 din stânga virgulei apărând întotdeauna, nu mai e necesar a fi memorat ceea ce permite creșterea cu o unitate a preciziei mantisei. El este însă implicit presupus, fiind restituit de coprocesor la începerea calculelor.

Semn	Exponent	Mantisă
------	----------	---------

Fig. 1.8 Formatul de reprezentare în virgulă flotantă

Exponentul este memorat sub formă cumulată. Valoarea de cumul (deplasamentul) a fiecărui format este o constantă și a fost astfel aleasă încât însumată cu exponentul să rezulte întotdeauna un exponent pozitiv. Aceasta simplifică foarte mult comparațiile. În final exponentul real se obține scăzând din exponentul numărului valoarea de cumul.

Tipurile de formate în virgulă flotantă sunt :

-real scurt (real în simplă precizie).

Acesta utilizează 32 biți din care unul pentru semn, 8 pentru exponent, 23 pentru mantisă. Mantisă e normalizată, iar bitul de parte întreagă este implicit 1, rezultând o precizie de 24 cifre binare. Valoarea de cumul este +127. Acest format permite memorarea unor numere cuprinse în intervalul $\pm 10^{\pm 38}$.

-real lung (real în dublă precizie).

Utilizează 64 biți din care unul este pentru semn, 11 pentru exponent, 52 pentru mantisă. Mantisă e normalizată, iar bitul de parte întreagă este 1 implicit, rezultând o precizie de 53 cifre binare. Valoarea de cumul este de +1023.

-real temporar (real de precizie extinsă).

Utilizează 80 biți din care: unul pentru semn, 15 pentru exponent, 64 pentru mantisă. Bitul de parte întreagă al mantisei normalizate este direct memorat, spre deosebire de celelalte două formate reale. Valoarea de cumul: + 16383. Este formatul utilizat intern de FPU pentru toate calculele. Eroarea de rotunjire și impreciziile de calcul se situează dincolo de limitele realului scurt și lung.

1.2.2 Realizarea celor patru operații aritmetice elementare în virgulă flotantă

Relațiile 1.11÷1.14 rezumă operațiile de bază ce se efectuează în aritmetica binară în virgulă flotantă. Pentru adunare și scădere trebuie asigurat ca cei doi operanzi să aibă același exponent. Aceasta implică deplasarea virgulei la unul dintre operanzi pentru a se realiza alinierea mantiselor. Înmulțirea și împărțirea sunt mult mai directe. Având astfel două numere în virgulă flotantă, X și Y:

$$X = M_x \times B^{E_x} \qquad Y = M_y \times B^{E_y}$$

atunci cele patru operații aritmetice elementare se efectuează în modul următor:

$$X+Y = (M_x \times B^{E_x-E_y} + M_y) \times B^{E_y} \qquad (1.11)$$

$$X - Y = (M_x \times B^{E_x-E_y} - M_y) \times B^{E_y}, \quad E_x \leq E_y \qquad (1.12)$$

$$X \times Y = (M_x \times M_y) \times B^{E_x+E_y} \qquad (1.13)$$

$$X / Y = (M_x / M_y) \times B^{E_x-E_y} \qquad (1.14)$$

în care B reprezintă baza de numerație în care se lucrează. În marea majoritate a aplicațiilor B=2.

Problemele care pot apărea în cazul acestor operații sunt:

- supradepășire exponent - un exponent pozitiv ar rezulta mai mare decât cea mai mare valoare posibilă a exponentului. În anumite sisteme acest caz poate duce la desemnarea lui ca $+\infty$ sau $-\infty$.
- subdepășire de exponent - un exponent negativ ar rezulta în modul mai mare decât cea mai mare valoare posibilă a exponentului. Aceasta înseamnă că numărul este prea mic pentru a putea fi reprezentat și poate fi interpretat ca fiind zero.
- subdepășire mantisă - în procesul de aliniere a mantiselor, biții se deplasează prea mult spre capătul din dreapta al mantisei, caz în care anumite metode de rotunjire sunt cerute.
- supradepășire mantisă - adunarea a două numere de același semn poate duce la obținerea unui transport peste cel mai semnificativ bit al mantisei. Acest lucru poate fi evitat prin realiniere.

Algoritmul de adunare și scădere în virgulă flotantă în manieră sincronă.

În aritmetica în virgulă flotantă adunarea și scăderea sunt mai complexe, din punct de vedere al etapelor ce trebuie parcurse, decât înmulțirea și împărțirea. Aceasta se datorează necesității de realizare a alinierii mantiselor pentru ca cei doi operanzi să prezinte același exponent. Algoritmul de adunare sau scădere cuprinde patru faze principale:

- Verificarea de zero
- Alinierea mantiselor
- Adunarea sau scăderea mantiselor
- Normalizarea rezultatului

Schema logică tipică ce realizează aceste funcții este dată în figura 1.9.

Metode de operare numerică

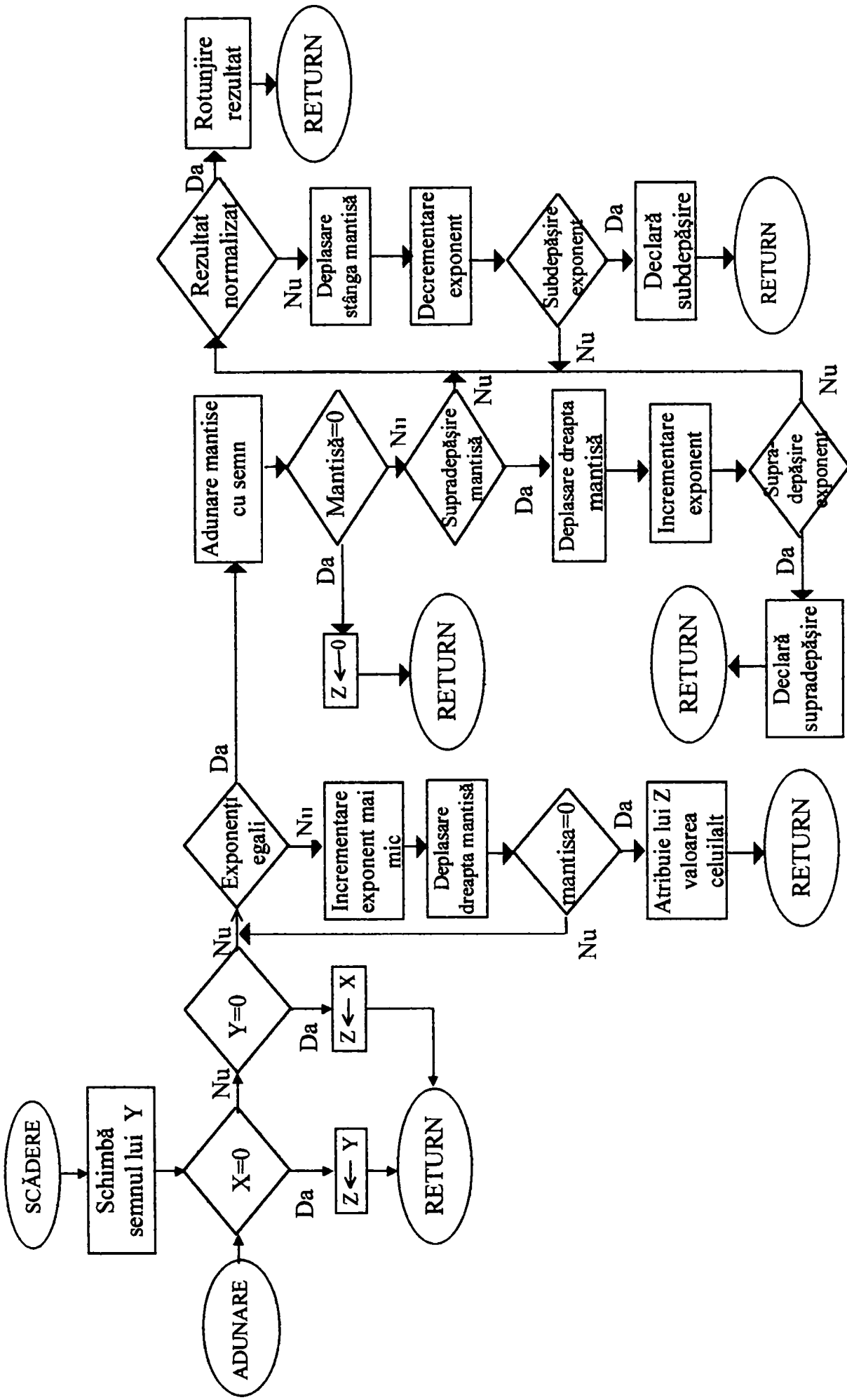


Fig. 1.9 Algoritm de adunare/scădere în virgulă flotantă

Pentru realizarea operațiilor de adunare sau scădere, cei doi operanzi trebuie transferați în regiștrii care vor fi utilizați de ALU. Dacă formatul în virgulă flotantă include implicit cel mai semnificativ bit, acesta trebuie să fie făcut explicit pentru realizarea operației. În mod obișnuit exponenții și mantisele vor fi memorați în regiștri separați, pentru a fi reușiți atunci când s-a produs rezultatul. Deoarece adunarea și scăderea sunt identice, cu excepția faptului că trebuie schimbat bitul de semn al scăzătorului în cazul scăderii, algoritmul este în rest același. Mai departe, dacă un operand este găsit zero, atunci celălalt operand este desemnat ca și rezultat. Următoarea fază este de a manipula astfel numerele încât exponenții să fie egali.

Alinierea mantiselor se realizează fie prin deplasarea celui mai mic număr spre dreapta, fie prin deplasarea celui mai mare număr spre stânga. Întrucât oricare din cele două variante ar putea duce la o pierdere de biți, cel mai mic număr se va deplasa. Orice bit care ar putea fi pierdut este astfel de relativ mică importanță. Alinierea este realizată prin deplasări repetate spre dreapta cu un bit și incrementarea exponentului până cei doi exponenți vor fi egali. (Dacă baza în care se lucrează este 16, deplasării cu un bit îi corespunde o deplasare cu câte 4 biți).

Dacă acest proces duce la obținerea unei mantise nule, celălalt număr este desemnat ca și rezultat. Astfel, dacă două numere prezintă exponenți care diferă semnificativ, atunci cel mai mic număr este pierdut.

În continuare cele două mantise sunt adunate împreună, luând în considerare și semnul lor. Din moment ce semnul poate diferi, rezultatul poate fi zero. De asemenea există posibilitatea ca mantisa rezultatului să producă o supradepășire cu un bit. Dacă e așa, mantisa se va deplasa spre dreapta cu un bit, iar exponentul se va incrementa o dată. Dacă exponentul produce supradepășire, operația se oprește.

Următoarea fază este normalizarea rezultatului. Normalizarea constă în deplasarea spre stânga a mantisei rezultatului până când cel mai semnificativ bit (sau 4 biți pentru exponent în baza 16) este nenul. Fiecare deplasare produce o decrementare a exponentului care ar putea duce la o subdepășire de exponent. În final rezultatul va fi rotunjit și reportat.

Algoritmul de înmulțire și împărțire în virgulă flotantă, în manieră sincronă.

Înmulțirea în virgulă flotantă, ca de altfel și împărțirea, sunt procese mult mai simple, din punct de vedere al numărului de etape ce trebuiesc parcurse, decât adunarea și scăderea. Operația de înmulțire este ilustrată în schema logică din figura 1.10.

La început, dacă fiecare operand este zero, rezultatul este declarat nul. Următorul pas este adunarea exponenților. Dacă exponenții sunt memorați în formă cumulată (cu deplasament), exponentul sumă va conține dublul deplasamentului, astfel că o valoare de deplasament trebuie scăzută din suma exponenților. Rezultatul ar putea produce fie supradepășire exponent, fie subdepășire exponent și acest lucru trebuie semnalizat la sfârșitul algoritmului. Dacă exponentul produsului este în domeniul bun, următorul pas constă în înmulțirea mantiselor, ținând cont și de semnul lor.

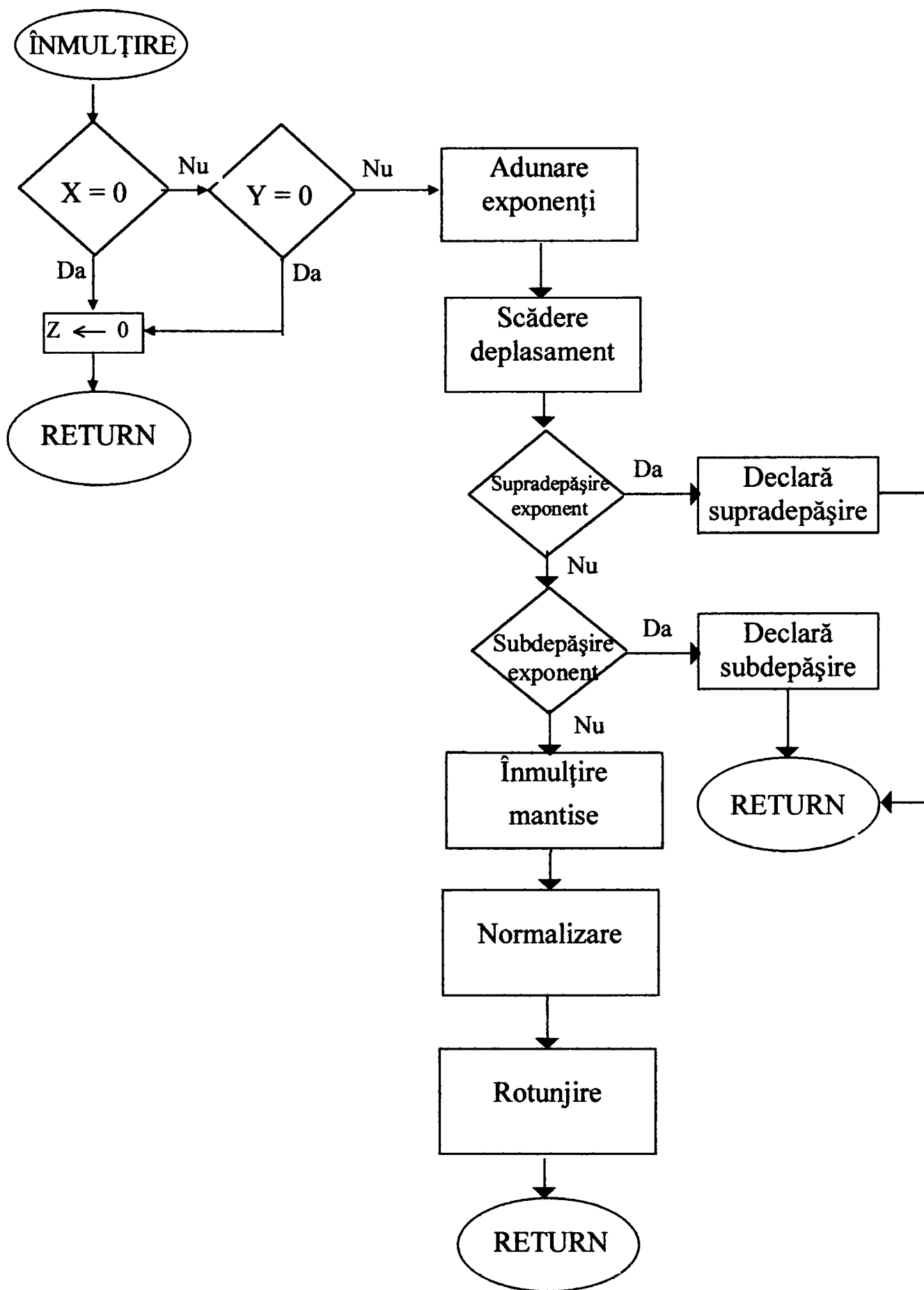


Fig. 1.10 Algoritmul de înmulțire în virgulă flotantă

Înmulțirea mantiselor este realizată în același mod ca și în cazul numerelor întregi. În acest caz avem de-a face cu o reprezentare semn - mărime, astfel că semnul rezultatului este stabilit prin intermediul unei porți SAU-exclusiv. Produsul obținut va avea o lungime egală cu suma lungimilor deînmulțitului și înmulțitorului. Biții rezultați peste lungimea maximă vor fi pierduți în timpul rotunjirii. Însumarea exponenților va ține cont de semnul fiecărui exponent.

După ce produsul este calculat, rezultatul este apoi normalizat și rotunjit așa cum s-a făcut și în cazul adunării, iar exponentul este ajustat, cu observația că normalizarea ar putea duce la o subdepășire de exponent.

Schema logică a algoritmului de împărțire în virgulă flotantă este dată în figura 1.11. Și aici primul pas constă în testarea operanzilor dacă au sau nu valoarea zero vreunul dintre ei. Dacă împărțitorul este zero este semnalizată eroare sau rezultatul este desemnat a fi infinit. Dacă deîmpărțitul este zero, atunci rezultatul este desemnat a fi zero.

În continuare exponentul împărțitorului este scăzut din exponentul deîmpărțitului. Aceasta va elimina valoarea de deplasament din exponent, astfel că deplasamentul trebuie adăugat din nou. De asemenea sunt făcute teste pentru subdepășirea sau supradepășirea exponentului.

Următorul pas este împărțirea mantiselor. Aceasta este apoi urmată de normalizare și rotunjire, respectiv ajustarea exponentului. Și în acest caz semnul rezultatului este stabilit prin intermediul unei porți SAU-exclusiv.

Considerații privind mecanismul de funcționare asincron/sincron al unei unități de calcul în virgulă flotantă.

Funcționarea în manieră pur sincronă a unei unități de calcul în virgulă flotantă, conform cu algoritmi prezentati anterior, conduce la efectuarea celor patru operații aritmetice elementare într-un timp foarte mare, întrucât ar fi necesare zeci sau chiar sute de perioade de tact pentru obținerea rezultatului. Realizarea pas cu pas a alinierii mantiselor sau a normalizării, în care deplasarea unui operand cu un bit este urmată de un test logic, conduce la o astfel de situație. Frecvența de clock a sistemului nu ar putea fi nici ea prea mare, fiind limitată de acea operație simplă în care informația străbate cele mai multe niveluri logice (ca de exemplu operația de adunare propriu-zisă a doi operanzi întregi).

Pe de altă parte, funcționarea în manieră pur asincronă a unei unități de calcul în virgulă flotantă ar conduce la realizarea oricărei operații aritmetice elementare într-o singură perioadă de tact. Timpul total de calcul ar rezulta mai scurt decât în cazul precedent, dar totuși mult mai lung decât în cazul unor operații uzuale de transfer de date între regiștrii coprocesorului, ceea ce pe ansamblu conduce la o anumită ineficiență precum și la viteză de operare scăzută.

Din acest motiv, o unitate de calcul în virgulă flotantă se organizează pe mai multe niveluri pipeline, separate de latch-uri. În interiorul unui nivel din structură informația se procesează în manieră asincronă, iar propagarea rezultatelor intermediare de la un nivel pipeline la altul se face sincron, în ritmul frecvenței de clock.

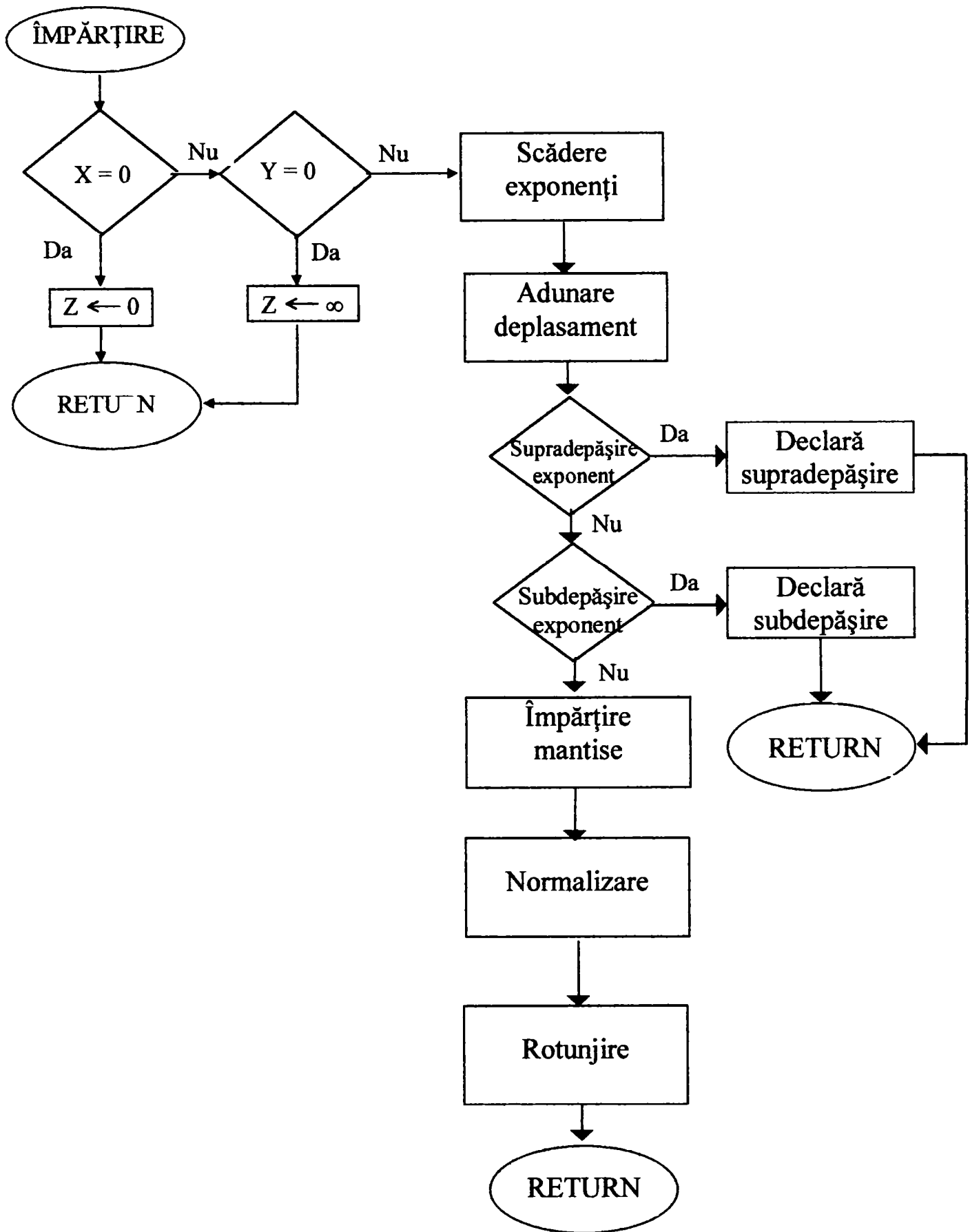


Fig. 1.11 Algoritm de împărțire în virgulă flotantă

1.2.3 Metode de creștere a vitezei operațiilor de adunare și scădere în virgulă flotantă

Întrucât multiplicatoarele binare în virgulă flotantă sunt mari consumatoare de arie pe cip, conform [14], [25], [29], [35], [55], [62], [65], [67], [70], [72], [81], [84], [97], [103], [106], iar circuitele de împărțire sunt lente comparativ cu circuitele care efectuează celelalte operații aritmetice elementare (de circa 4÷6 ori mai lente decât multiplicatoarele binare) conform [24], [53], [54], [66], [69], [73], [80], ele vor fi implementate folosind principiile aritmeticii logaritmice, așa cum se va vedea în capitolul următor. În schimb operarea pe baza circuitelor de adunare/scădere în virgulă flotantă reprezintă modalitatea de găsire a rezultatului operațiilor amintite cea mai rapidă și mai economică din punct de vedere al ariei ocupate pe cip. Din acest motiv, în acest paragraf vor fi amintite principalele soluții descrise în literatura de specialitate, care concură la implementarea unor circuite de adunare/scădere în virgulă flotantă performante.

Așa cum s-a remarcat în paragraful anterior, structura optimă a unei unități de calcul în virgulă flotantă presupune o organizare pe niveluri pipeline a acesteia. O proiectare judicioasă a unității de calcul trebuie să conducă la situația în care toate nivelurile pipeline să fie parcurse de informația binară în intervale de timp aproximativ egale și acestea să fie cât mai scurte. În plus structura pipeline permite, în cazul operării consecutive a mai multor perechi de operanzi, obținerea unui nou rezultat la fiecare perioadă de clock a sistemului.

În cazul circuitelor de adunare/scădere în virgulă flotantă, soluția clasică, încă foarte frecvent adoptată, presupune o organizare pe trei niveluri pipeline a acestora (**soluția I**), dar tendința actuală este de a se reduce numărul de niveluri la două (**soluția II**) în condițiile păstrării sau chiar creșterii frecvenței de clock a sistemului. Prețul plătit constă în aproape dublarea numărului de porți logice integrate, fapt care însă este suportat tehnologic datorită creșterii densității de integrare.

Etapele logice succesive care sunt parcurse la o adunare sau scădere în virgulă flotantă sunt următoarele:

- a. Calcularea diferenței exponenților;
- b. Stabilirea exponentului mai mic;
- c. Deplasarea corespunzătoare a mantisei asociate exponentului mai mic;
- d. Operarea (adunarea/scăderea) mantiselor având asociate exponentul mai mare;
- e. Determinarea numărului de zerouri după virgulă conținute în diferența mantiselor, respectiv a rangului celui mai semnificativ 1 din partea întreagă a sumei mantiselor;
- f. Normalizarea rezultatului;
- g. Ajustarea exponentului rezultatului;
- h. Rotunjirea rezultatului.

La elaborarea **soluției I** s-a avut în vedere în principal reducerea numărului de niveluri logice, prin introducerea unor structuri de calcul paralele și tehnici de

selecție a rezultatului corect în interiorul fiecărui nivel pipeline. Au putut fi astfel suprapuse aproape în totalitate, din punct de vedere temporal, anumite etape din lista enunțată care în mod natural se desfășoară succesiv; imediat după stabilirea definitivă a rezultatului de la finele etapei anterioare s-a procedat la selectarea rezultatului corect.

Astfel, în **nivelul 1** al structurii pipeline, scăderea exponenților (etapa a.) se face concomitent cu deplasarea în paralel a celor două mantise (etapa c.), ținând cont că circuitele “barrel shifter”, care deplasează în manieră asincronă mantisele, sunt comandate pe măsură ce se produc biții diferenței exponenților, dinspre LSB spre MSB, conform [26], [66]. În momentul în care se cunoaște deja semnul diferenței exponenților (etapa b.), se selectează mantisa deplasată ce a corespuns exponentului mai mic.

De asemenea, în **nivelul 2** al structurii pipeline, operația de scădere a mantiselor (etapa d.), prin maniera în care a fost implementată, elimină timpul necesar reconversiei din cod complement de doi în cod semn-mărime, în cazul diferenței negative, conform [7], [26], [39], [44], [56]. Acest lucru se realizează prin calcularea simultană a două rezultate, din care doar unul va reprezenta valoarea corectă și va fi selectat în funcție de semnul diferenței celor două mantise. De asemenea operația de rotunjire (etapa h.) prin adaos, sau către cea mai apropiată valoare reprezentabilă în format, se poate efectua în acest nivel, prin calcularea simultană a încă unui rezultat.

Atunci însă când se efectuează scădere propriu-zisă și operanzii sunt foarte apropiați în valoare absolută, pot rezulta multe zerouri înaintea celui mai semnificativ 1 din rezultat. Pentru a se efectua normalizarea (etapa f.), aceste zerouri trebuie în primă fază contorizate (etapa e.), urmând ca apoi să fie comandat un circuit “barrel shifter” similar cu cel de la alinierea mantiselor.

O economie importantă de timp se face dacă numărul acestor zerouri (“leading zeros”) poate fi anticipat încă din etapa de operare a celor două mantise aliniate, astfel ca la finele acestei etape să înceapă efectiv deplasarea în vederea normalizării. O serie de articole și lucrări științifice din literatura de specialitate au avut ca temă această problemă [7], [10], [85]. Soluțiile oferite au permis, prin adăugarea unor structuri hard suplimentare în nivelul 2 al structurii pipeline, să se estimeze numărul de zerouri cu o precizie de 1 bit, ceea ce face ca, în urma comandării circuitului “barrel shifter” de normalizare în **nivelul 3**, rezultatul să se găsească în domeniul $[0.5, 2)$. Cu alte cuvinte, se obține fie rezultatul normalizat, fie având un zero înaintea celui mai semnificativ 1. În urma testării bitului MSB al acestuia se mai comandă un circuit suplimentar de postnormalizare, care fie lasă rezultatul nedeplasat, fie execută o deplasare spre stânga cu un bit. În felul acesta, pe ansamblu, s-a făcut o economie de timp semnificativă. Tot în acest nivel se efectuează în paralel și operația de ajustare a exponentului rezultatului (etapa g.).

Totuși operația de postnormalizare nu este chiar foarte scurtă datorită necesității multiplicării capabilității de comandă a bitului MSB menționat, ce trebuie să comande un număr foarte mare de grile MOS, din multiplexoarele 2:1

ale circuitului de postnormalizare. De aceea, au fost concepute și circuite de corecție a erorii de predicție a numărului de zerouri [11], care elimină astfel timpul consumat la postnormalizare. Evident, acest lucru presupune și o creștere suplimentară semnificativă a numărului de porți logice utilizate.

La elaborarea **soluției II** s-a avut în vedere reducerea numărului de niveluri logice, prin introducerea unor structuri de calcul paralele și tehnici de selecție a rezultatului corect nu doar în interiorul fiecărui nivel pipeline ci la scara întregii unități de adunare/scădere în virgulă flotantă. Rezultă astfel două căi distincte de procesare a informației, care nu se reunesc decât în urma selecției finale, de la ieșirea acestei unități de calcul. Acest mod de abordare pleacă de la următoarele premize:

- în cazul operației de adunare, normalizarea ulterioară a rezultatului este foarte scurtă. Într-adevăr, având mantisa asociată exponentului mai mare, normalizată în intervalul $[1,2)$, deoarece bitul de rang zero este 1 implicit, iar cealaltă mantisă, deplasată sau nedeplasată, e situată în intervalul $[0,2)$, rezultă că suma nu se poate găsi decât în intervalul $[1,4)$. Aceasta înseamnă că după testarea bitului de rang 1 din partea întregă a rezultatului, se decide dacă suma rămâne nedeplasată sau se deplasează la dreapta cu un bit, în vederea normalizării;

- în cazul operației de scădere, dacă operanzii sunt mult diferiți ca și valoare, este necesar un timp important pentru deplasarea mantisei mai mici, dar rezultatul operării mantiselor aliniate se va situa în intervalul $[0.5, 2)$ și normalizarea va fi foarte scurtă. În acest caz fie se lasă rezultatul nedeplasat, fie se deplasează la stânga cu un bit. Dacă în schimb operanzii sunt foarte apropiați, alinierea mantiselor este foarte rapidă, în sensul că fie se deplasează mantisa mai mică cu un bit la dreapta, fie se lasă mantisele nedeplasate. În acest caz operația de normalizare este cea mai lungă, întrucât rezultă multe zerouri înaintea celui mai semnificativ 1 din diferența mantiselor. Pe ansamblu, timpul de efectuare a operației de scădere este aproximativ același în cele două cazuri.

Concluzia care se desprinde constă în aceea că etapa de aliniere a exponenților și deplasare corespunzătoare a mantisei mai mici și respectiv etapa de normalizare sunt complementare din punct de vedere al timpului necesar pentru parcurgerea lor, în cazul cel mai defavorabil al fiecăreia dintre ele. Astfel, în ultimele lucrări apărute în literatura de specialitate [7], [77] se definesc și se implementează două căi distincte ("far path" sau cale îndepărtată și "near path" sau cale apropiată) de operare pentru cele două cazuri prezentate, ceea ce presupune practic aproape dublarea ariei circuitelor de adunare/scădere în virgulă flotantă. Sunt necesare chiar și două sumatoare/scăzătoare distincte deoarece în unitatea de calcul în virgulă flotantă nivelul la care intervin aceste circuite pe cele două căi este diferit. Ele pot fi chiar secționare între două etaje consecutive, în scopul optimizării vitezei, ceea ce presupune și utilizarea unor latch-uri suplimentare. În aceste condiții, circuitul de adunare/scădere în virgulă flotantă poate fi organizat pe doar două niveluri pipeline, fără a reduce frecvența de clock a sistemului.

1.3 Metode de operare numerică în aritmetică logaritmică

1.3.1 Realizarea operațiilor de înmulțire și împărțire a două numere cu ajutorul logaritmilor

Operațiile de înmulțire și împărțire, efectuate cu ajutorul logaritmilor, sunt bazate pe următoarele relații:

$$A \times B = \text{antilog} (\log A + \log B) \quad (1.15)$$

$$A / B = \text{antilog} (\log A - \log B) \quad (1.16)$$

în care A și B reprezintă modulul celor 2 numere care se înmulțesc. Semnul produsului este generat în exteriorul structurii hard care permite implementarea relațiilor precedente.

Pentru a păstra o gamă dinamică largă de reprezentare a numerelor, cei doi factori A și B vor fi reprezentați în format virgulă flotantă cu mantisa normalizată :

$$A = x \times 2^i \quad B = y \times 2^j \quad (1.17)$$

unde $x \geq 1$, $y \geq 1$, $x, y < 2$. În mod uzual este utilizată baza 2 pentru reprezentarea numerelor, iar logaritmul va fi de asemenea considerat în bază 2.

Logaritmii produsului și ai câtului sunt dați de relațiile (1.18 și 1.19):

$$\log (A \times B) = i + j + \log xy \quad (1.18)$$

$$\log (A / B) = i - j + \log x/y \quad (1.19)$$

iar produsul, respectiv câtul sunt obținute pe baza relațiilor (1.20) și (1.21):

$$A \times B = 2^{i+j} \text{antilog} (\log x + \log y) = 2^{i+j} \times z \quad (1.20)$$

$$A / B = 2^{i-j} \text{antilog} (\log x - \log y) = 2^{i-j} \times z' \quad (1.21)$$

Astfel generarea produsului (câtului) actual necesită o operație de adunare (scădere) realizată atât pentru logaritmii mantiselor cât și pentru exponenți, implementată cu sumatoare (scăzătoare) rapide, două operații de logaritmare ce se pot realiza simultan și o operație de antilogaritmare.

Logaritmare și antilogaritmare pot fi realizate fie prin utilizarea unor memorii ROM care memorează direct valoarea logaritmilor și antilogaritmilor binari conform [9] (**metoda 1**), fie prin utilizarea unor structuri logice celulare (SLC) de generare a logaritmilor și antilogaritmilor binari conform [16], [17] (**metoda 2**), fie utilizând algoritmi rapizi ce necesită memorii ROM de mult mai mică capacitate, care memorează valoarea logaritmului, respectiv antilogaritmului în anumite puncte între care se face interpolare liniară [40], [41], [42], [44], [50], [51], [107] (**metoda 3**).

Principalii indicatori de performanță pentru un circuit de înmulțire logaritmice sunt următorii:

- eroarea cu care se calculează logaritmul și antilogaritmul, respectiv produsul final;

Metode de operare numerică

- timpul sau viteza de generare a logaritmului și antilogaritmului, respectiv a produsului final;
- aria ocupată pe cipul de siliciu;
- gradul de modularitate al circuitului, sau ușurința implementării în layout.

Acești indicatori vor fi analizați pe scurt în cazul celor trei metode amintite.

Metoda 1, descrisă în [9], este cea mai simplă și mai directă. Eroarea care apare la generarea lui z sau z' din relațiile (1.20) și (1.21) este aceeași deoarece adunarea și scăderea produc aceeași eroare absolută și se determină, doar pentru cazul înmulțirii, pe baza următoarelor considerente :

- prima sursă de eroare este dată de însăși reprezentarea numărului x pe n biți, obținându-se $x_{(n)}$:

$$e_1 = x_{(n)} - x \quad (1.22)$$

- a doua sursă de eroare apare la generarea pe m biți a logaritmului binar $\log_{(m)} x$:

$$e_2 = \log_{(m)} x - \log x \quad (1.23)$$

- a treia sursă de eroare apare la generarea antilogaritmului binar pe n biți (s-a considerat că produsul este reprezentat pe n biți, ca și factorii):

$$e_3 = \text{antilog}(\log x + e_2) + e_1 - x = \text{antilog}(\log_{(m)} x) + e_1 - x \quad (1.24)$$

Deoarece e_2 provine dintr-o reprezentare pe m biți a unui număr subunitar ($0 < \log x < 1$), rezultă că $e_2 \leq \frac{1}{2} \frac{1}{2^m}$. Deoarece e_3 provine dintr-o reprezentare pe n biți a unui număr cuprins între 1 și 4, (produsul a două numere cuprinse între 1 și 2), rezultă că $e_3 \leq \frac{1}{2} \frac{4}{2^n} = \frac{1}{2} \frac{1}{2^{n-2}}$. În aceste condiții, eroarea totală care apare la executarea operației de înmulțire după această metodă este:

$$\begin{aligned} e &\leq \text{antilog}(\log x_{(n)} + \log y_{(n)} + \frac{1}{2} \frac{1}{2^m} + \frac{1}{2} \frac{1}{2^m}) + \frac{1}{2} \frac{1}{2^{n-2}} - x_{(n)} y_{(n)} = \\ &= \text{antilog}(\log x_{(n)} y_{(n)} + \frac{1}{2^m}) + \frac{1}{2^{n-1}} - x_{(n)} y_{(n)} = x_{(n)} y_{(n)} [\text{antilog}(\frac{1}{2^m}) - 1] + \frac{1}{2^{n-1}} \end{aligned} \quad (1.25)$$

Modelul matematic complet al înmulțirii logaritmice poate fi reprezentat astfel în figura 1.12:

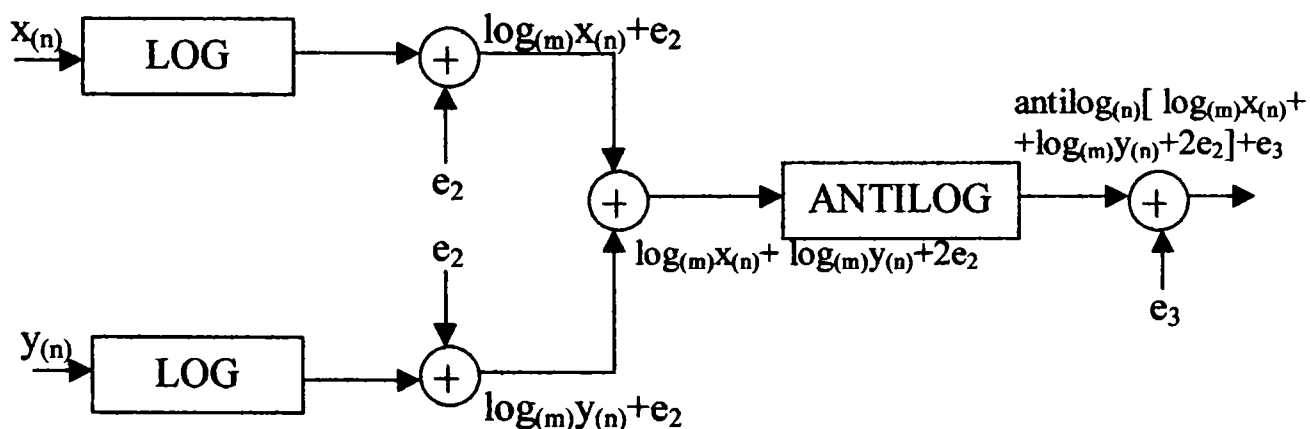


Fig. 1.12 Modelul matematic complet al înmulțirii logaritmice

Analizând relația (1.25), se constată că eroarea totală absolută depinde de valoarea operanzilor și va avea valoarea maximă atunci când operanzii sunt foarte aproape de valoarea 2. Pentru acest caz se dezvoltă în serie Taylor expresia erorii totale, obținându-se:

$$e \leq \frac{1}{2^{n-1}} + \left(\frac{\ln 2}{2^{m-2}} + \frac{1}{2} \frac{(\ln 2)^2}{2^{2m-2}} + \dots \right) \cong \frac{1}{2^{n-1}} + \frac{\ln 2}{2^{m-2}} \quad (1.26)$$

Se observă din relația (1.26) că atunci când $m \rightarrow \infty$, deci când logaritmiile operanzilor sunt calculați și reprezentați cu precizie infinitesimală, eroarea totală rămâne cea din cazul când se efectua înmulțirea propriu-zisă a doi operanzi reprezentați pe n biți. Pentru a estima numărul de biți cu care trebuie reprezentat logaritmul pentru ca eroarea totală să rămână acceptabilă, a fost reprezentată grafic în figura 1.13 eroarea totală pentru cazurile când $n=22$, $n=23$ și $n=24$, în funcție de m .

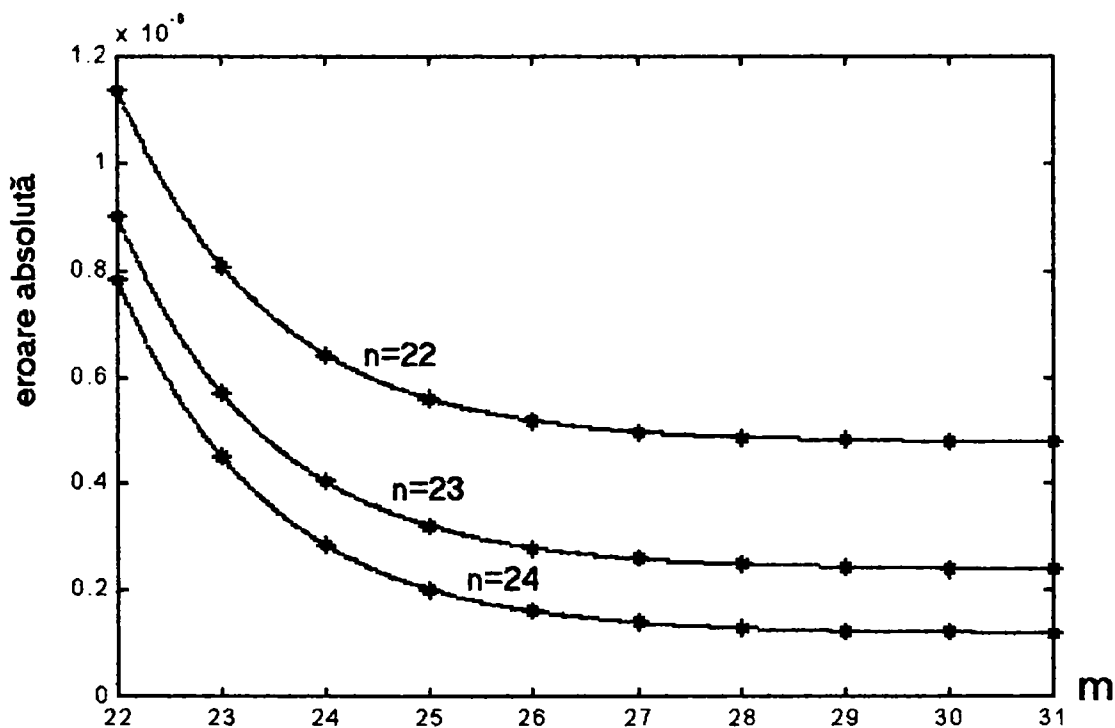


Fig. 1.13 Eroarea totală absolută în cazul înmulțirii logaritmice

Diagramele de genul celor din figura 1.13 sunt foarte utile la stabilirea formatului operanzilor și la dimensionarea capacității memoriilor atunci când este impusă o anumită valoare maximă a erorii de înmulțire. De exemplu, în cazul în care se lucrează în simplă precizie ($n=23$), pentru ca eroarea totală să fie aproximativ aceeași cu cea de la înmulțirea propriu-zisă, numărul de biți pe care trebuie reprezentat logaritmul este 27 sau chiar 28.

Inconvenientul major pe care îl prezintă **metoda 1** constă în cantitatea uriașă de memorie care trebuie integrată. Volumul total de memorie necesară, în biți, este $2 \times m \times 2^n + n \times 2^m$, atunci când sunt utilizate 2 memorii ROM pentru generarea simultană a logaritmilor lui x și y . Tabelul 1.2 prezentat mai jos, împreună cu figura 1.13 permit determinarea dimensiunii optime a memoriei ROM, exprimată în Mbit, din diversele combinații ale lungimii cuvântului produsului și

logaritmului. Așa cum se observă din acest tabel, este necesară o structură multicip pentru a se putea integra cantitatea de memorie necesară, ceea ce conduce și la pierderea, într-o anumită măsură, a avantajului dat de viteza superioară de generare a logaritmului și antilogaritmului prin citirea directă din ROM.

Tab. 1.2

$m \rightarrow$	22	23	24	25	26	27	28	
$n \downarrow$	22	277 Mb	377 Mb	570 Mb	948 Mb	1694 Mb	3179 Mb	6140 Mb
	23	466 Mb	579 Mb	789 Mb	1191 Mb	1980 Mb	3540 Mb	6644 Mb
	24	839 Mb	973 Mb	1208 Mb	1644 Mb	2483 Mb	4127 Mb	7382 Mb

Metoda 2, prezentată în [16] și [17], este bazată pe descompunerea numărului ce se logaritmează în factori de forma $(1+k_i 2^{-i})$, respectiv pe descompunerea numărului ce se antilogaritmează în termeni de forma $k_i L_i$. Coeficienții k_i pot fi 0 sau 1 și se numesc coeficienți de creștere, ei arătând dacă un factor respectiv termen este prezent sau nu în descompunere. Valorile L_i reprezintă logaritmii binari, ai numerelor de forma $1+2^{-i}$. Valoarea maximă a lui i determină precizia cu care se calculează logaritmul sau antilogaritmul.

Structura hard care implementează această metodă prezintă următoarele caracteristici:

- operează cu numere normalizate, fiind compatibilă cu sistemele de calcul în virgulă flotantă;
- schema sa are o structură celulară și deci poate fi ușor reprodusă și extinsă;
- durata maximă de generare a logaritmului sau antilogaritmului este aproximativ $500 \times t_p$ (t_p - timpul de propagare printr-o poartă logică simplă) în simplă precizie;
- eroarea apărută la logaritmare sau antilogaritmare se poate face oricât de mică prin simpla creștere a numărului de celule.

Avantajul adus prin folosirea acestei metode constă în aceea că aria circuitelor crește doar cu pătratul numărului de biți ai operanzilor și nu exponențial, ca în cazul metodei precedente. Un dezavantaj important este însă dat de timpul prea mare de propagare a transportului prin structură, care nu poate fi micșorat semnificativ folosind sumatoare rapide, datorită algoritmului iterativ folosit, pe bază de comparații repetate.

Metoda 3 de calcul al logaritmului și antilogaritmului binar are la bază un algoritm ce va fi prezentat pe larg în capitolul 2 al tezei și care prin contribuțiile aduse de autor, permite o viteză de operare în formatul simplă precizie de circa 10 ori mai mare decât în cazul metodei precedente. Aria ocupată pe cip este redusă, iar eroarea furnizată este comparabilă cu a sistemelor în virgulă flotantă.

1.3.2 Realizarea unor operații matematice complexe cu ajutorul logaritmilor binari.

Reprezentarea numerelor binare în format logaritmnic (LNS - Logarithm Number System) este o reprezentare în virgulă fixă având n biți, din care i biți pentru partea întreagă și f biți pentru partea fracționară.

Proprietățile logaritmilor binari permit calculul foarte simplu pentru o serie de funcții matematice care în virgulă flotantă necesită un timp mult mai îndelungat întrucât se apelează fie la algoritmi specifici complecși, fie la dezvoltări în serie care pentru precizia dorită ajung să conțină foarte mulți termeni. În continuare se va arăta pe scurt cum anumite funcții matematice pot fi calculate foarte simplu folosind aritmetica logaritmică.

-calculul rădăcinii pătrate. Dacă x este un număr pozitiv atunci:

$$\sqrt{x} = \text{antilog}(\log \sqrt{x}) = \text{antilog}\left(\frac{1}{2} \log x\right) \quad (1.27)$$

Se știe că împărțirea la 2, în binar, echivalează cu o deplasare la dreapta cu un bit. În felul acesta se pierde cel mai puțin semnificativ bit iar în locul celui mai semnificativ bit se introduce un zero. Astfel, calculul rădăcinii pătrate se obține printr-o logaritmare, o deplasare la dreapta cu un bit și o antilogaritmare. Similar, calculul rădăcinii de ordin 4 implică o deplasare a logaritmului numărului x cu doi biți la dreapta. Calculul rădăcinii de ordin 2^y (y fiind un întreg pozitiv) implică o deplasare a logaritmului numărului x cu y biți la dreapta. În felul acesta se pierd cei mai puțin semnificativi y biți iar în locul celor mai semnificativi y biți se introduc zerouri. Întrucât astfel ar avea loc o anumită pierdere de precizie, reprezentarea LNS se face în interiorul procesorului logaritmic pe un număr ceva mai mare de biți decât la interfața cu procesorul principal.

-calculul pătratului.

$$x^2 = \text{antilog}(\log x^2) = \text{antilog}(2 \log x) \quad (1.28)$$

Cu alte cuvinte calculul pătratului se face printr-o logaritmare, o deplasare la stânga cu un bit și o antilogaritmare. Similar pentru calculul unei puteri, putere de 2 (2^y), se va face deplasarea spre stânga cu y biți.

- calculul oricărei puteri cu exponent real (x , R pozitive).

$$x^{\pm R} = \text{antilog}(\log x^{\pm R}) = \text{antilog} \{ \pm \text{antilog}[\log(\log x) + \log r] \} \quad (1.29)$$

Așadar calculul acestei funcții complexe se face în doar 5 pași: o dublă logaritmare pentru x și $\log x$ simultan cu $\log R$, o adunare, urmată apoi de două antilogaritmări consecutive.

- calculul unui logaritm în orice bază.

$$\log_B x = \frac{\log x}{\log B} = \text{antilog} [\log(\log x) - \log(\log B)] \quad (1.30)$$

Acest calcul se face așadar prin două logaritmări consecutive simultane, pentru x și $\log x$, respectiv pentru B și $\log B$, o scădere și o antilogaritmare, adică în doar 4 pași.

Se observă că toate aceste funcții matematice fac în cele din urmă apel la operații de logaritmare și antilogaritmare pentru care s-au conceput, așa cum se va vedea mai târziu, algoritmi competitivi și la adunare de numere întregi, pentru care au fost prezentate deja sumatoare rapide performante. Operația de scădere se reduce tot la adunare folosind reprezentarea în complement de doi a scăzătorului.

1.3.3 Realizarea operațiilor de adunare și scădere în aritmetica logaritmică

Dificultatea principală în realizarea unei unități aritmetice logaritmice complete constă în implementarea operațiilor de adunare și scădere în special. Implementările anterioare anului 1990, [52], [89], au fost limitate la un număr redus de biți de precizie (sub dimensiunea formatului simplă precizie) datorită faptului că aria circuitelor utilizate în aritmetica logaritmică crește exponențial cu lungimea cuvintelor de date vehiculate, atunci când nu se folosesc algoritmi de interpolare. Ulterior, progresul tehnologic (creșterea densității de integrare) dar și conceperea unor algoritmi pentru calculul logaritmului și antilogaritmului cât și pentru calculul funcțiilor logaritmice specifice necesare realizării operațiilor de adunare și scădere au permis implementarea unor unități logaritmice complete în simplă precizie [57], [58], [59], [104]. Într-o tehnologie submicronică, aceste unități de calcul chiar pot fi integrate într-un singur cip.

Calculul funcțiilor menționate anterior este făcut utilizând interpolarea liniară între anumite valori memorate ale acestor funcții, facilitând o reducere drastică (de peste 100 de ori) a necesarului de memorie ROM, în condițiile în care eroarea de calcul este încă acceptabilă, comparativ cu sistemele de calcul în virgulă flotantă.

În continuare va fi prezentat pe scurt algoritmul prin care se implementează operațiile de adunare și scădere în aritmetica logaritmică.

Considerăm un număr x care este reprezentat de perechea (S_x, e_x) în care S_x dă informația cu privire la semnul lui x , iar e_x este logaritmul său. Astfel $x = (-1)^{S_x} \times 2^{e_x}$ unde e_x este un număr reprezentat în virgulă fixă având N biți, cu I biți pentru partea întreagă și F biți pentru partea fracționară. $N = F + I$. Valoarea lui e_x este dată de expresia (1.31):

$$e_x = \sum_{i=-F}^{I-1} e_{xi} \times 2^i \quad (1.31)$$

Pentru a obține reprezentarea lui 0, se utilizează un bit distinct Z_x ca și indicator de zero. Dacă $Z_x = 1$ atunci valoarea numărului este zero. Astfel, numărul x este reprezentat în formatul LNS de tripleta (S_x, Z_x, e_x) și are valoarea :

$$x = (1 - Z_x) \times (-1)^{S_x} \times 2^{e_x} \quad (1.32)$$

Formatul LNS simplă precizie ce include bit de semn și bit indicator de zero presupune în consecință $N = 30$, cu $F = 22$ și $I = 8$.

Dacă operațiile de înmulțire sau împărțire se rezumă la adunarea sau scăderea exponenților, operațiile de adunare și scădere sunt mai complexe. Considerăm A și B două numere pozitive și $A \geq B$. Nu se pierde din generalitate ci doar se simplifică notațiile.

În cazul adunării ($C = A + B$), logaritmul e_c al sumei C are valoarea:

$$e_c = e_a + f_a(e_a - e_b) = e_a + f_a(r), \text{ în care: } f_a(r) = \log(1 + 2^{-r}) \quad (1.33)$$

În cazul scăderii ($C = A - B$), logaritmul e_c al diferenței C are valoarea:

$$e_c = e_a + f_s(e_a - e_b) = e_a + f_s(r), \text{ în care: } f_s(r) = \log(1 - 2^{-r}) \quad (1.34)$$

Funcția “log” reprezintă logaritmul în baza 2.

Metoda clasică pentru implementarea funcțiilor $f_a(r)$ și $f_s(r)$ constă în utilizarea unor memorii ROM având valorile posibile ale lui r ca și intrări și furnizând la ieșire cuvintele $f(r)$. Într-o implementare directă, pentru formatul simplă precizie descris anterior, ar rezulta însă un număr de cuvinte $2^N = 2^{30}$ care reprezintă o cantitate uriașă de memorie, imposibil de realizat practic. Reducerea semnificativă a necesarului de memorie se face pe baza formei particulare de variație a funcțiilor $f_a(r)$ și $f_s(r)$ reprezentate în figura 1.14.

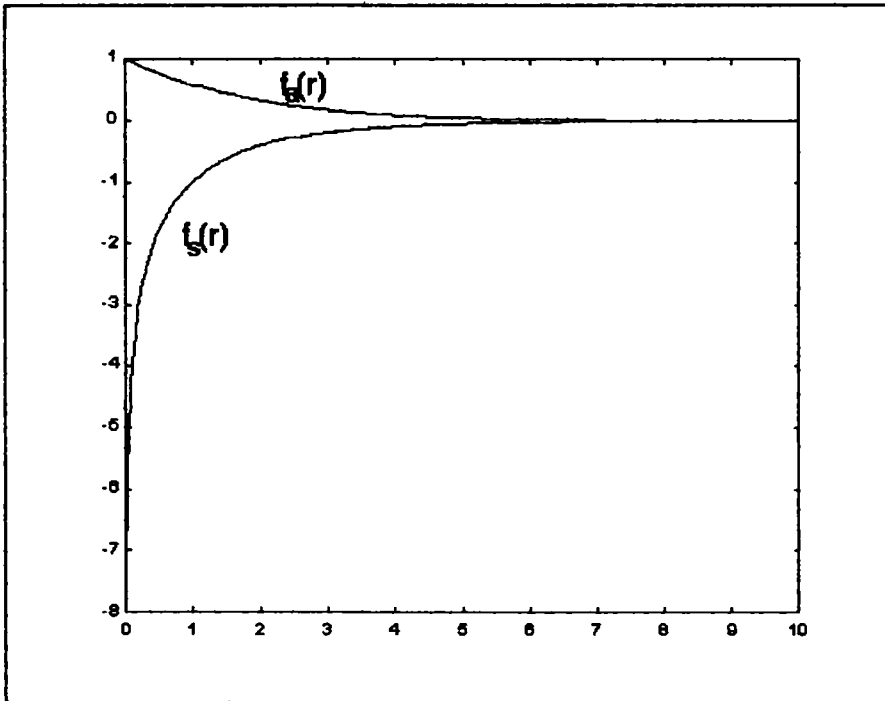


Fig. 1.14 Graficul funcțiilor $f_a(r)$ și $f_s(r)$

Atunci când argumentul r (care e un număr pozitiv) depășește o anumită limită, valorile funcțiilor $f_a(r)$ și $f_s(r)$ sunt atât de aproape de zero încât valoarea lor nu mai poate fi reprezentabilă (coboară sub rezoluția sistemului). Astfel, dacă $r = F+2$ atunci rezultă:

$$\begin{aligned} |f_a(r)| &< 2^{-F-1}, \\ |f_s(r)| &< 2^{-F-1}, \end{aligned}$$

în care 2^{-F-1} este precizia cu care poate fi reprezentat un număr în format LNS.

Astfel memoria ROM trebuie implementată doar pentru valori ale lui r care satisfac dubla inegalitate : $0 < r \leq F+2$ ceea ce în cazul concret cu $F=22$ înseamnă $0 < r \leq 24$. De asemenea se demonstrează că pentru o precizie dată cu F biți pentru partea fracționară, în cazul unui r care satisface dubla inegalitate $i \leq r < i+1$, cu $i+1 \leq 24$, ultimii cei mai puțin semnificativi i biți nu pot afecta rezultatul furnizat de tabela funcției f implementată în memoria ROM. Acest lucru permite partiționarea domeniului util al lui r în subintervale egale cu unitatea și implementarea pentru fiecare funcție (f_a și f_s) a F tabele separate (pentru fiecare subinterval) conținând între 1 și 2^F cuvinte, ceea ce în total înseamnă aproximativ 4×2^F cuvinte în loc de 2^N . Aceasta extinde posibilitatea implementării prin memorare directă până la maxim 16÷18 biți de precizie fracționară. Utilizând însă aproximarea seriei Taylor pentru funcțiile f_a și f_s , se poate extinde numărul biților de precizie fracționară la valoarea $F=22$.

Considerând funcția $f : \mathbb{R} \rightarrow \mathbb{R}$, diferențiabilă, aproximarea Taylor de ordin n a acestei funcții, f^* , în vecinătatea lui r este :

$$f^*(r + \Delta r) = f(r) + \sum_{i=1}^n \frac{(\Delta r)^i}{i!} \times \frac{d^i f(r)}{dr^i} \quad (1.35)$$

Aproximarea Taylor liniară care va fi folosită utilizează numai primii doi termeni din expresia (1.35). Argumentul r al funcției se consideră acum discretizat, fiind reprezentat pe N biți și se efectuează o partiție a cuvântului r în patru câmpuri distincte după cum urmează :

$$r = r_i + r_h + r_l + r_e \quad (1.36)$$

r_i conține partea întreagă a lui r deci conține I biți: $r_i = r_{I-1} \dots \dots \dots r_0$;

r_h conține " p_l " biți: $r_h = r_{-1} \dots \dots \dots r_{-p_l}$;

r_l conține " $(p_e - p_l)$ " biți: $r_l = r_{-p_l-1} \dots \dots \dots r_{-p_e}$;

r_e conține " $F - (p_e)$ " biți: $r_e = r_{-p_e-1} \dots \dots \dots r_{-F}$.

În continuare, se definește $r_t = r_i + r_h$.

Valorile lui p_l și p_e nu vor fi constante . Ele vor fi alese în funcție de subintervalul în care se găsește r . Aproximarea Taylor va utiliza ca și argument r pe r_t iar ca și cantitate Δr pe r_l . Cantitatea r_e va fi ignorată, în conformitate cu observația făcută la analiza necesarului de memorie. Astfel se poate scrie :

$$f^*(r) = f(r_t) + r_l \times \frac{df(r_t)}{dr} \quad (1.37)$$

Memoria ar conține deci două tabele, una pentru valorile lui f , alta pentru cele ale lui df/dr care ar conține 2^{p_l} , respectiv $2^{p_e - p_l}$ cuvinte. Totuși relația (1.37) presupune utilizarea unui multiplicator. Evitarea acestuia se face apelând încă o dată la aritmetica logaritmică. Notând $2^x = \exp(x)$ se poate rescrie relația (1.37):

$$f^*(r) = f(r_t) + \text{sgn}\left(\frac{df(r_t)}{dr}\right) \times \exp\left(\log\left(\left|\frac{df(r_t)}{dr}\right|\right) + \log(r_l)\right) \quad (1.38)$$

Proprietățile particulare ale derivatelor funcțiilor f_a și f_s permit simplificarea în continuare a relației (1.38). Astfel:

$$\frac{df_a}{dr} = \frac{2^r}{1+2^r}; \quad \log\left|\frac{df_a}{dr}\right| = r - \log(1+2^r) = r - f_a(r) \quad (1.39)$$

$$\text{Similar se obține: } \log\left|\frac{df_s}{dr}\right| = r - f_s(r) \quad (1.40)$$

Astfel, se obțin expresiile (1.41) și (1.42) care indică modalitatea de implementare a operațiilor de adunare și scădere în format LNS:

- în cazul adunării:

$$e_c = e_a + f_a(r_t) + \exp(\log(r_l) + r_t - f_a(r_t)) \quad (1.41)$$

- în cazul scăderii:

$$e_c = e_a + f_s(r_t) - \exp(\log(r_l) + r_t - f_s(r_t)) \quad (1.42)$$

Algoritmul prin care se realizează aceste operații este prezentat în tabelul 1.3 în care s-au notat și explicitat rezultatele operațiilor intermediare:

Tab. 1.3

adunare	scădere
$f_{rt} = f_a.tbl(r_t)$	$f_{rt} = f_s.tbl(r_t)$
$l_{rl} = \log.tbl(r_l)$	$l_{rl} = \log.tbl(r_l)$
$f_d = r_t - f_{rt}$	$f_d = r_t - f_{rt}$
$l_{cor} = l_{rl} + f_d$	$l_{cor} = l_{rl} + f_d$
$cor = \exp.tbl(l_{cor})$	$cor = \exp.tbl(l_{cor})$
$ff = f_{rt} + cor$	$ff = f_{rt} - cor$
$e_c = e_a + ff$	$e_c = e_a + ff$

S-a ținut cont de faptul că pentru $\forall r > 0$, $df_a/dr > 0$ și $df_s/dr < 0$.

Arhitectura structurii hard ce realizează operația de adunare și scădere în aritmetică logaritmică este dată în figura 1.15.

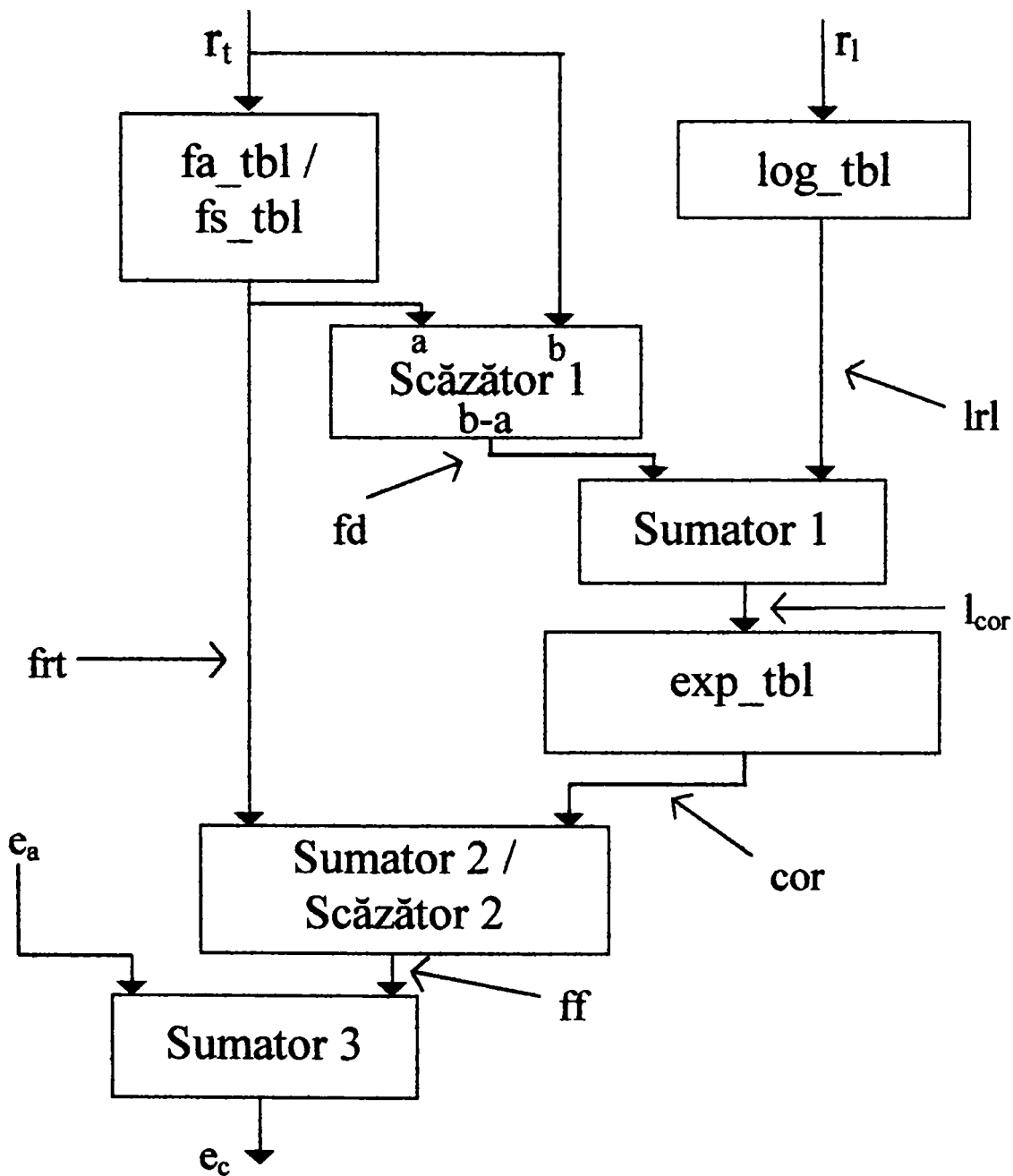


Fig.1.15 Schema bloc a circuitului de adunare/scădere LNS

Capacitatea memoriilor care furnizează pe $\log(r_1)$ și $\exp(l_{cor})$ este considerabil mai mică decât a acelor care furnizează valorile funcțiilor f_a și f_s . Maniera de programare a memoriilor fa_tbl și fs_tbl nu este totuși perfect asemănătoare datorită faptului că funcția f_s prezintă o singularitate pentru $r=0$. Din acest motiv intervalul $0 < r \leq 1$ a fost subîmpărțit și el în mai multe intervale. Valorile lui p_1 și p_e depind de valoarea lui r relativ la aceste intervale.

Două tipuri de erori pot afecta precizia operațiilor de adunare și scădere :

- a) eroarea datorată utilizării aproximării Taylor de ordinul 1 și care este aproximativ egală cu diferența dintre aproximarea Taylor de ordinul 2 și aproximarea f^* utilizată pentru funcțiile f ;
- b) erorile care apar în fiecare etaj atunci când datele parcurg structura prezentată în figura 1.15 și care generează o eroare cumulată la ieșirea întregului bloc. Această eroare apare datorită faptului că fiecare tabelă de memorie utilizată are o mărime finită și că fiecare cuvânt vehiculat pe calea de date este reprezentat pe un număr finit de biți.

Unele dintre aceste erori contribuie direct la eroarea finală, iar o altă parte, inclusiv din prima categorie contribuie și indirect întrucât ele vor trece prin tabela care furnizează pe 2^x (exp_tbl).

În lucrarea [58] este prezentată metoda de calcul pentru optimizarea diverselor lungimi de cuvinte și pentru deducerea valorilor p_e și p_1 , în condițiile în care eroarea totală este 2^{-F-1} , adică cea corespunzătoare formatului simplă precizie. Metoda a constatat în alocarea fiecărui termen din expresia dedusă pentru eroarea totală, a câte unei părți din disponibilul de " 2^{-F-1} ", alocare efectuată pe baza următoarelor priorități:

-e preferabilă micșorarea lungimii unui cuvânt care reprezintă o adresă de intrare în memorie decât a lungimii unui cuvânt ce reprezintă o ieșire din memorie. Aceasta deoarece un bit adițional la intrarea de adresă dublează capacitatea memoriei pe când un bit adițional în cuvântul de ieșire îi mărește capacitatea doar cu un număr de biți egal cu numărul de cuvinte conținute;

-în reducerea dimensiunii memoriilor, prioritate vor avea memoriile care furnizează pe $f_a(r)$ și $f_s(r)$ deoarece sunt componentele hard cele mai expansive. Aceasta implică selectarea în primul rând a termenilor din expresia erorii totale care îl conțin pe p_1 ce determină direct numărul de intrări în tabelele ROM menționate; următoarele componente hard semnificative sunt tabelele ROM log_tbl și exp_tbl .

Capacitatea memoriilor care furnizează pe f_a și f_s a rezultat $(F+10) \times 2^{F/2} - 8$ cuvinte de $F+7$ biți, iar capacitatea totală de memorie $(F+18) \times 2^{F/2} - 8$ cuvinte de diferite lungimi, rezultând pentru $F=22$, o capacitate de 2,66 Mbit. Această cantitate de memorie poate fi integrată într-o tehnologie submicronică într-un singur cip, împreună cu celelalte componente ale procesorului aritmetic logaritmic.

Viteza de operare a circuitului de adunare/scădere logaritmic este însă de $2 \div 3$ ori mai redusă decât a unui circuit omolog în virgulă flotantă. Așa cum se

observă din figura 1.15, calea critică de propagare este cea pe care se succed blocurile: fa_tbl/fs_tbl , Scăzător 1, Sumator 1, exp_tbl , Sumator 2/Scăzător 2 și Sumator 3. În tehnologie CMOS 0,5 micrometri, doar timpul de acces pentru fiecare din cele două memorii de pe calea critică de propagare este de peste 5ns [107]. Utilizând un mecanism "CSA-carry save adder" (vezi paragraful 2.2.2.1 și figura 2.5) și metoda implementată în circuitul din figura 2.24 pentru a comasa blocurile Scăzător 1 și Sumator 1, respectiv Sumator 2/Scăzător 2 și Sumator 3, pot fi obținuți timpi de operare pentru blocurile comasate de ordinul a 4 ns. În aceste condiții, timpul total de generare a rezultatului este de minim 18 ns ($5+4+5+4$), în condițiile în care o adunare/scădere în virgulă flotantă se poate efectua în $6\div 8$ ns [7].

În plus, atât operanzii de intrare cât și rezultatul se găsesc în format LNS, deci la timpul de calcul estimat se mai adaugă și timpul necesar conversiei din și înapoi în formatul virgulă flotantă, dacă aplicația în care este utilizat procesorul aritmetic necesită acest lucru. Chiar și din punct de vedere al ariei ocupate pe cip, raportul este net favorabil sistemului de operare în virgulă flotantă.

Totuși, datorită simplității prin care se realizează operațiile de înmulțire și împărțire și deci a vitezei obținute în cazul operării în format LNS, un procesor de calcul pur logaritmic merită a fi propus în anumite aplicații particulare cum ar fi:

- implementarea unor algoritmi complecși, în care rezultatele intermediare pot rămâne în format LNS și numai un operand trebuie adus la intrarea sau de la ieșirea sistemului;

- aplicații în care nu este necesar a se adopta un format standardizat și deci nu este necesară nici o conversie inițială de format întrucât mărimea analogică ce este achiziționată poate fi furnizată de un sistem cu caracteristică de transfer logaritmică;

- în sistemele de calcul complexe (multiprocesor), în care conversia de format poate fi realizată în timpul transferului DMA în bufferele de intrare sau ieșire.

Cap.2 Concepția subunității logaritmice a unui procesor aritmetic hibrid, ce permite efectuarea operațiilor de înmulțire și împărțire

În acest capitol, precum și în cel următor, va fi prezentată și aplicată o nouă concepție privind proiectarea unui procesor aritmetic care va îmbina proprietățile aritmeticii logaritmice și avantajele operării în virgulă fixă cu calitățile reprezentării numerelor în virgulă flotantă și de operare cu acestea folosind cei mai performanți algoritmi de calcul.

2.1 Premize teoretice. Deziderate propuse

Așa cum s-a precizat în Introducere, tipul unui coprocesor este dat de modul în care sunt reprezentate datele în interiorul său. Atunci când simultan sunt cerute atât o scară largă de reprezentare a numerelor reale cât și precizie, în mod uzual se adoptă sistemul de reprezentare în virgulă flotantă, motiv pentru care procesoarele aritmetice care operează cu date de acest tip sunt cele mai răspândite.

Creșterea, în plan tehnologic, a densității de integrare a permis dezvoltarea, ca și alternativă și a procesoarelor logaritmice, ce operează cu date în format LNS. Au apărut astfel o serie de unități de calcul logaritmice din ce în ce mai performante, fiind menționate și în literatura de specialitate: [2], [4], [12], [18], [21], [22], [37], [40], [42], [44], [45], [50], [51], [57], [58], [59], [89], [90], [104]. Singurul lor dezavantaj este că deocamdată nu permit decât calcule în simplă precizie. Din acest motiv și reprezentarea în virgulă flotantă nu va face apel în capitolul următor decât la formatul simplă precizie.

Un procesor aritmetic hibrid, în virgulă flotantă și logaritmice (FP-LNS) este un procesor care vehiculează date în interiorul său, reprezentate în ambele tipuri de formate, pe căi specifice fiecareia dintre ele. Prof. M. Ciugudean enunță acest concept în 1976 în teza de doctorat, [16], în care erau oferite soluțiile pentru efectuarea logaritmirii și antilogaritmirii și în care era prezentată arhitectura unui procesor hibrid și modul de implementare a diverselor operații matematice. Ulterior, (1986) a reluat și a dezvoltat această problematică [15], [18].

Prima implementare pe cip a unui astfel de procesor (1985) e datorată lui F. J. Taylor [90]. Unitatea de execuție numerică (NEU) a acestui coprocesor realiza operațiile matematice de bază în format LNS. Legătura dintre NEU și unitatea de control (CU) a coprocesorului, ce vehicula date în virgulă flotantă era făcută prin intermediul a două convertoare de format: FP-LNS (pentru logaritmare) și LNS-FP (pentru antilogaritmare). Dezavantajul principal care se făcea resimțit era dat de dimensiunea mică a cuvintelor de date cu care se opera. Întrucât conversiile de format utilizau memorii ROM interne care furnizau valoarea logaritmului sau antilogaritmului prin citire directă, capacitatea acestora nu putea fi foarte mare, ceea ce limita lungimea cuvintelor de date la o valoare inferioară standardului existent în FP. În același timp, lungimea cuvintelor de date era restricționată și de capacitatea memoriilor ROM ce trebuiau integrate pentru implementarea

operațiilor de adunare și scădere în format logaritmice. Ulterior, Lewis prin algoritmul propus în 1990 [58], beneficiind și de evoluția tehnologică, oferă o soluție pentru realizarea operațiilor de adunare și scădere LNS în simplă precizie, dar nu este făcută nici o specificație cu privire la realizarea logaritmării respectiv antilogaritmării.

Ultimele mențiuni în literatura de specialitate cu privire la concepția unui procesor aritmetic hibrid sunt făcute de Lai în 1991 [50], [51] și 1993 [107]. Lai a propus și un algoritm avantajos pentru efectuarea conversiilor de format. *Acest algoritm va fi adoptat și în lucrarea de față cu unele modificări avantajoase, efectuate de autorul tezei, menite să reducă complexitatea logicii de comandă a ALU din subunitatea logaritmice de calcul precum și să elimine o parte din circuitele de prelucrare a exponenților din convertoarele de format, asigurând și o viteză de calcul ceva mai mare. Algoritmul va fi implementat, tot ca și contribuție proprie, printr-o proiectare completă până la nivel de poartă logică.*

În continuare, pe baza analizei efectuate asupra algoritmului cât și asupra întregii structuri proiectate după acest algoritm, vor fi aduse, atât la nivel de principiu cât și hard, contribuții suplimentare importante. Va fi descrisă și implementată o nouă metodă de structurare a arhitecturii unității logaritmice. Implementarea se va realiza utilizând circuite de concepție proprie, mai performante decât cele descrise în literatura de specialitate. De asemenea programarea memoriilor ROM utilizate la generarea logaritmului și antilogaritmului va utiliza o metodă suplimentară de corecție a erorilor cauzate de metoda de interpolare cunoscută în bibliografie.

Toate acestea, așa cum se va dovedi prin analiză Matlab și simulare, vor conduce atât la o creștere de precizie cât și la o creștere substanțială a vitezei de calcul a logaritmului și respectiv antilogaritmului, în condițiile în care numărul de porți logice utilizate, respectiv numărul de tranzistoare integrate pe cip este și mai redus.

Concepția procesorului aritmetic hibrid ce va fi prezentat în acest capitol a plecat de la următoarele premize:

-procesoarele logaritmice sunt foarte eficiente în operații ca înmulțirea, împărțirea, ridicarea la putere, extragerea de radical, calculul funcției logaritm în orice bază sau al funcției de ridicare la orice putere reală. În schimb complexitatea operațiilor de adunare și scădere realizate în format LNS, așa cum s-a văzut în paragraful 1.3, conduce la o pierdere inevitabilă de viteză și la consum excesiv de arie pe cipul de siliciu datorită capacității mari de memorie ROM ce trebuie integrată.

-procesoarele în virgulă flotantă implementează algoritmul de adunare/scădere în format FP, prezentat în paragraful 1.2 al tezei, prin intermediul unor structuri hard cu funcționare asincronă, astfel că, în ciuda complexității sale comparativ cu aritmetica în virgulă fixă, rezultatul poate fi furnizat într-un timp relativ scurt. Înmulțirea, de asemenea, se realizează foarte rapid (cu ceva mai mult timp decât adunarea), tot în manieră asincronă, utilizând așa zisele arii de înmulțire

[80]. Aceasta însă se face cu prețul unui consum mare de arie pe cip [30], [35], [55], [62], [65], [67], [72], [81], [97], [103]. În ce privește împărțirea, aceasta este atât consumatoare de suprafață cât și de circa 4-6 ori mai lentă decât înmulțirea [50], [53], [69], [73], întrucât stabilirea biților cântului este întretesută cu operații de comparare/scădere. Calculul pentru toate celelalte funcții matematice, inclusiv rădăcina pătrată, logaritmul și exponențiala, se face mult mai lent decât în cazul adunării și scăderii în virgulă flotantă.

În concluzie, procesorul hibrid descris în acest capitol va realiza adunarea și scăderea în virgulă flotantă, în simplă precizie, (suficientă încă în multe aplicații) și toate celelalte operații prin intermediul unei subunități logaritmice de calcul. Aceasta conduce la o arhitectură originală ce presupune că însăși unitatea de execuție numerică (NEU) a procesorului aritmetic va avea o structură hibridă, iar convertoarele de format FP-LNS și LNS-FP se vor situa în interiorul său.

Pentru creșterea vitezei de calcul și asigurarea posibilității implementării cu ușurință a unor algoritmi DSP, subunitățile coprocesorului hibrid au fost proiectate în structură pipeline.

Întrucât atât subunitatea în virgulă flotantă cât și cea logaritmice a NEU efectuează operații care în cele din urmă se reduc la adunare de numere întregi, vor fi adoptate sumatoare binare performante, cunoscute în literatură, însă *pe căile critice de propagare vor fi utilizate sumatoare de concepție proprie, mai rapide.*

În plus, s-a căutat suprapunerea pe cât posibil a unor operații care în mod natural se desfășoară succesiv, respectiv combinarea implementării unor structuri binare paralele cu tehnici de selecție a rezultatului corect. Astfel, în cazul adunării în virgulă flotantă, scăderea exponenților se va face în paralel cu deplasarea mantisei, ținând cont că circuitul care deplasează mantisa ("barrel shifter") este comandat pe măsură ce se produc biții diferenței exponenților, dinspre LSB spre MSB, conform [66]. De asemenea, operația de scădere a mantiselor, prin maniera în care a fost implementată, elimină timpul necesar reconversiei din cod complement de doi în cod semn-mărime, în cazul diferenței negative, conform metodei descrise în [26], dar *uzând de o topologie îmbunătățită a circuitelor utilizate, concepută de autor. De această proprietate va beneficia și circuitul de adunare/scădere din secțiunea logaritmice, prin intermediul căruia se vor realiza operațiile de înmulțire/împărțire și care a fost optimizat în urma modificării efectuate de autor a algoritmului propus de Lai în [50].*

2.2 Subunitatea de calcul logaritmice pentru realizarea rapidă a operațiilor de înmulțire și împărțire în simplă precizie

2.2.1 Algoritm de realizare a conversiilor de format FP-LNS și LNS-FP

Algoritm de realizare a conversiilor de format FP-LNS și LNS-FP a fost descris de către F. Lai în [50] și [51] și indică de fapt maniera prin care se poate efectua logaritmic și respectiv antilogaritmic.

În formatul virgulă flotantă (FP) simplă precizie un număr binar A , este reprezentat astfel:

S_A -1 bit	E_A -8 biți	M_A -23 biți
--------------	---------------	----------------

în care S_A este bitul de semn, E_A este exponentul, iar M_A mantisa normalizată. Deoarece întodeauna există un singur 1 în stânga virgulei, acesta nu se mai memorează dar este implicit restituit în interiorul coprocesorului în structurile hard care operează cu mantise. Valoarea numărului A este dată de următoarea expresie:

$$A = (-1)^{S_A} (1 + 0, M_A) 2^{E_A - 127} \quad (2.1)$$

În formatul simplă precizie, conform standardului IEEE754, valoarea de deplasament este 127 și se adaugă la exponentul real al numărului (formând exponentul cumulat), pentru a avea întotdeauna un exponent pozitiv.

În format logaritmice (LNS) în simplă precizie, un număr binar z este reprezentat astfel:

S_z -1 bit	I -8 biți	F -23 biți
--------------	-------------	--------------

Aceasta este o reprezentare în virgulă fixă cu $i=8$ biți pentru partea întreagă I , cu $f=23$ biți pentru partea fracționară F și un bit de semn S_z . Valoarea numărului z este dată de următoarea expresie:

$$z = (-1)^{S_z} 2^{N_z}, \quad N_z = I + F, \quad (2.2)$$

în care N_z este logaritmul binar al numărului pozitiv $|z|$. În mod uzual se lucrează cu logaritmul în baza 2 deoarece este cel mai ușor de obținut, [16] și în plus majoritatea procesoarelor aritmetice operează cu numere în baza 2.

Considerând mantisa din formatul în virgulă flotantă ca fiind normalizată, adică în intervalul $[1,2)$, partea întreagă a logaritmului numărului este dată chiar de valoarea exponentului fără deplasament, iar partea fracționară de logaritmul mantisei. Această concatenare este posibilă deoarece logaritmul mantisei este întotdeauna un număr subunitar pozitiv. Într-adevăr, rescriind relația 2.1 după modelul relației 2.2 se obține:

$$A = (-1)^{S_A} 2^{E_A - 127 + \log(1 + 0, M_A)}, \quad N_z = E_A - 127 + \log(1 + 0, M_A) \quad (2.3)$$

și identificând termenii lui N_z rezultă: $I = E_A - 127$ și $F = \log(1 + 0, M_A)$.

Datorită lungimii cuvintelor de date, nu a putut fi utilizată metoda simplă a citirii din memorie a logaritmului sau antilogaritmului fiecărui număr reprezentabil în formatul în simplă precizie întrucât ar fi fost nevoie de integrarea unei cantități enorme de memorie. Din acest motiv s-a utilizat o tabelă implementată într-o memorie ROM care va furniza doar anumite valori ce vor reprezenta valori de corecție în metoda utilizată și între care se va realiza interpolare liniară. Calculul logaritmului și antilogaritmului utilizează aproximarea lui Mitchell, folosită și de Maenner în [63], însă mult îmbunătățită.

a) Conversia FP-LNS

Se efectuează partiția argumentului și memorarea doar a anumitor valori pentru reducerea cantității de memorie așa cum a sugerat Maenner [63] și în plus se aplică o metodă de corecție care se bazează pe memorarea în anumite puncte și a valorilor derivatei funcției, după care se realizează interpolarea.

Considerând numărul A dat de relația 2.1 și făcând notația $0, M_A = y$, se poate scrie :

$$N_Z = \log_2 |A| = E_A - 127 + \log_2(1+y) \quad (2.4)$$

Semnul se procesează separat, exponentului i se scade valoarea de cumul (deplasamentul de 127) deci conversia care mai trebuie făcută este de fapt calculul funcției $\log_2(1+y)$ pe 23 de biți, în care $0 \leq y < 1$.

Se partiționează mantisa y în două cantități : y_1 conținând cei mai semnificativi 11 biți și y_2 conținând cei mai puțin semnificativi 12 biți. Curba erorii de conversie în cazul utilizării aproximării lui Mitchell, arătată în figura 2.1.a, se discretizează în $2^{11} = 2048$ puncte ale căror proiecții pe abscisă vor fi echidistante. Valorile în aceste puncte ale funcției $\log_2(1+y)-y$ sunt memorate într-un ROM ca și valori de corecție E_y , furnizate prin aplicarea adresei y_1 .

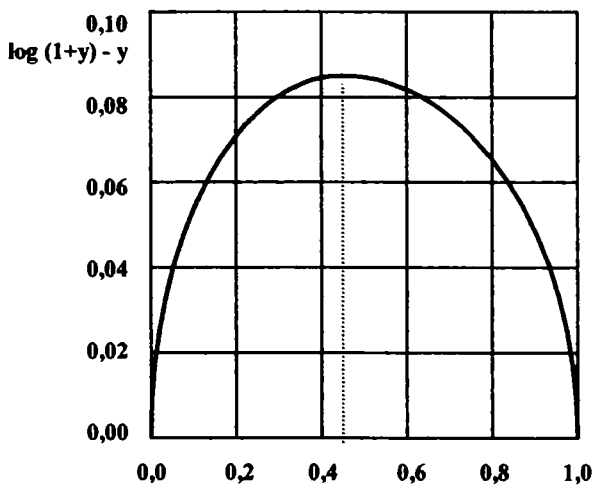


Fig.2.1.a

Curba erorii de conversie la logaritmare, în cazul utilizării aproximării lui Mitchell

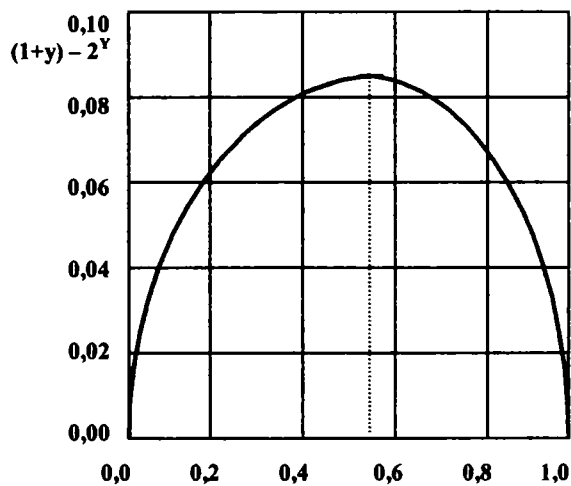


Fig.2.1.b

Curba erorii de conversie la antilogaritmare, în cazul utilizării aproximării lui Mitchell

Următoarea aproximare, îmbunătățită, care se poate face acum, este :

$$\log_2(1+y) \approx y + E_y \quad (2.5)$$

Totuși în afară de cele 2048 puncte care furnizează valoarea exactă a funcției (afectată doar de precizia reprezentării pe 23 de biți) celelalte valori rămân încă afectate de eroare. Reducerea în continuare a erorilor se face prin aplicarea interpolării liniare între puncte consecutive memorate. Pentru aceasta se memorează în ROM și diferențele ΔE_y dintre valorile învecinate ale funcției $\log_2(1+y)-y$. Obținem astfel următoarea aproximare :

$$\log_2(1+y) \approx y + E_y \pm \Delta E_y \times y_2 \quad (2.6)$$

A doua tabelă ROM este mult mai mică întrucât $\Delta E_y \ll E_y$. Considerând pentru ΔE_y o reprezentare pe 12 biți, conversia completă între cele două formate se face printr-o citire în tabelele de memorie, două adunări pe 24 de biți și o înmulțire de 12×12 biți. Vor fi utilizate sumatoare rapide și un multiplicator utilizând sumatoare de tip CSA aranjate în arbore Wallace.

b. Conversia LNS-FP

Procesul de antilogaritmare se tratează similar. Considerând x rezultatul antilogaritmării, se poate scrie:

$$x = 2^{(E-127)+(0.M)} = 2^{(E-127)} 2^{0.M} = 2^{(E-127)} 2^y \quad (2.7)$$

în care $E-127$ reprezintă partea întreagă în format LNS, iar M reprezintă partea fracționară. Se face următoarea aproximare:

$$x = 2^{(E-127)} 2^y \approx 2^{(E-127)}(1+y) \quad (2.8)$$

Se observă că și în acest caz mantisa în format FP este, în primă instanță, același lucru cu partea fracționară M din format LNS. Cifra 1 este bitul implicit din formatul FP. Această primă aproximare dă o eroare de conversie de valoare $(1+y) - 2^y$ a cărei curbă de variație este dată în figura 2.1.b; y se partiționează în același mod și se utilizează o memorie ROM pentru a memora eroarea de conversie E_y în 2048 puncte precum și diferența ΔE_y . Rezultatul final al conversiei este :

$$2^y \approx (1+y) - E_y \pm \Delta E_y \times y_2 \quad (2.9)$$

Prin utilizarea algoritmilor pentru conversiile de format FP-LNS și LNS-FP, operații ca înmulțirea sau împărțirea a două numere A și B pot fi reduse la adunări sau scăderi în virgulă fixă conform relațiilor de mai jos:

$$C = A \times B = \text{antilog}(\log A + \log B) \quad (2.10)$$

$$C = A / B = \text{antilog}(\log A - \log B) \quad (2.11)$$

În cazul numerelor cu semn, acesta se procesează separat utilizând o poartă SAU-exclusiv:

$$S_C = S_A \oplus S_B \quad (2.12)$$

Plecând de la relațiile 2.6, 2.9, 2.10 și 2.11, F. Lai definește arhitectura unei unități de calcul logaritmice [50], concepută pentru efectuarea operațiilor de înmulțire (MUL) și împărțire (DIV) în cazul a doi operanzi de intrare. Această arhitectură este prezentată în figura 2.2 și ilustrează maniera în care este prelucrată informația de la intrările unității de calcul. În figura 2.2, A' și B' corespund reprezentării în format FP cu exponent cumulat a celor două numere A și B , y_A și y_B sunt mantisele lor normalizate, y_{1A} și y_{1B} sunt cantitățile cele mai semnificative din y_A și y_B de lungime 11 biți, y_{2A} și y_{2B} sunt cantitățile cele mai puțin semnificative din y_A și y_B de lungime 12 biți, E_A și E_B sunt exponenții cu

deplasament, S_A și S_B sunt biții de semn, F_A și F_B sunt părțile fracționare ale operanzilor în format LNS, I_A și I_B sunt părțile lor întregi, I_C și F_C sunt partea întreagă respectiv fracționară, în format LNS a rezultatului furnizat de ALU, iar E_C , y_C și S_C sunt respectiv exponentul cu deplasament, mantisa și bitul de semn al rezultatului înmulțirii/împărțirii celor doi operanzi de intrare, în format FP cu exponent cumulat.

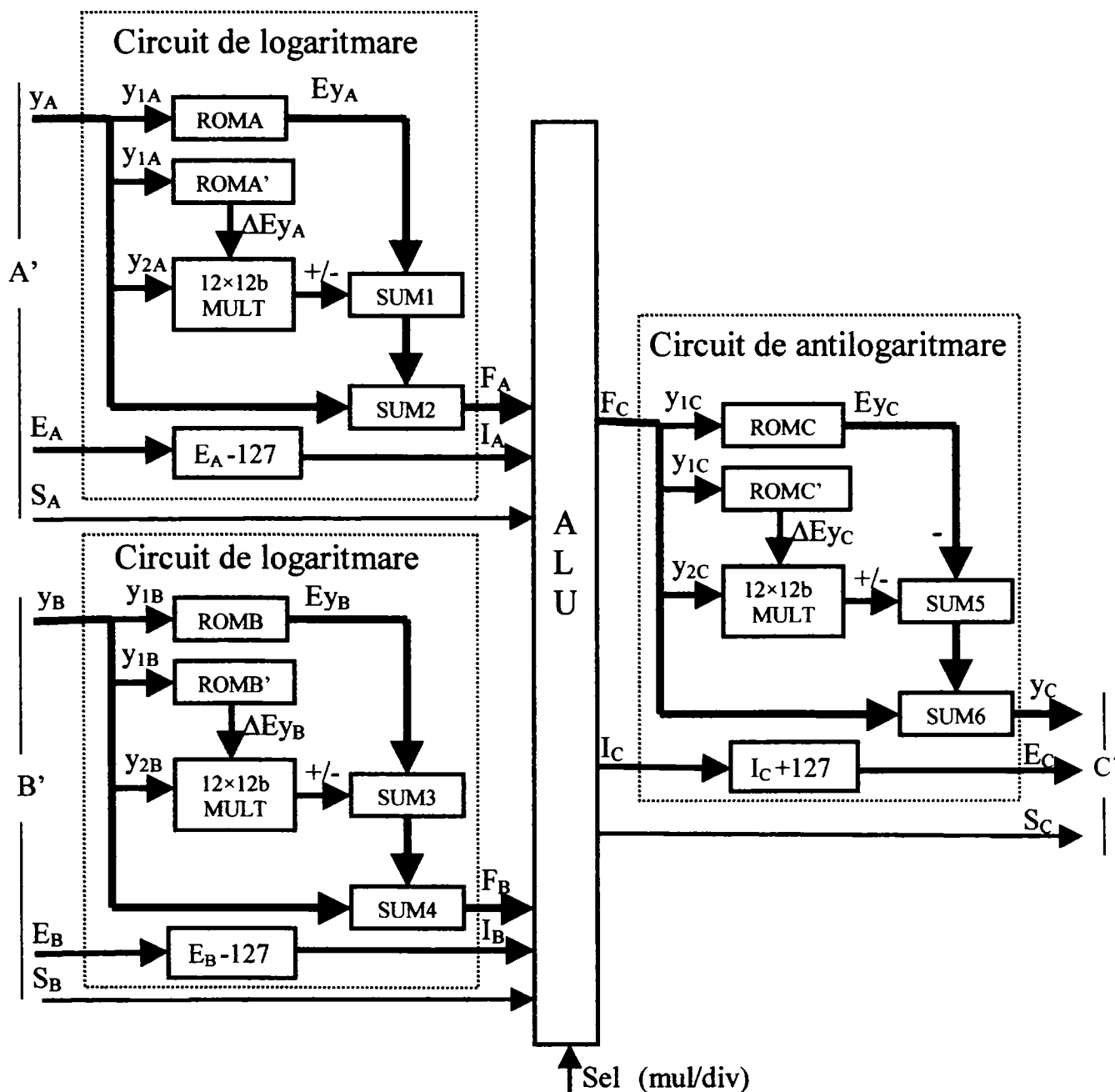


Fig. 2.2. Arhitectura subunității logaritmice ce implementează algoritmul lui F.Lai pentru realizarea conversiilor de format FP-LNS și LNS-FP

Memoriile ROMA și ROMA', accesate cu același cuvânt de adresă, sunt identice cu ROMB și respectiv ROMB', accesate în aceeași manieră. Așa cum se vede în figură, sunt necesare 2 memorii de tip ROMA+ROMA' și o memorie de tip ROMC+ROMC'. Blocul ALU nu conține altceva decât o poartă SAU-exclusiv

pentru implementarea relației 2.12 și un circuit sumator/scăzător pe 31 de biți pentru implementarea relației:

$$I_C, F_C = I_A, F_A \pm I_B, F_B \quad (2.13)$$

în care operația executată depinde de valoarea bitului Sel.

Arhitectura subunității de calcul logaritmice propusă de F. Lai în [50] pentru realizarea operațiilor de înmulțire și împărțire, așa cum se observă în figura 2.2, presupune utilizarea a trei sumatoare pe 8 biți, pentru a extrage din exponenții celor 2 operanzi deplasamentul de 127, respectiv pentru a adăuga valoarea de deplasament la exponentul rezultatului. Deoarece în acest fel se vehiculează exponenți care pot fi și negativi, rezultă că ALU va opera cu date în virgulă fixă ce pot fi de orice polaritate. Așadar apare în mod obligatoriu necesitatea asigurării conversiei în cod complement de doi pentru oricare din cele două date. În plus, atunci când se execută operația de împărțire care presupune scădere în ALU, trebuie făcut apel încă o dată la conversia în cod complement de doi a scăzătorului. Iar dacă rezultatul se dovedește a fi negativ, el este automat reprezentat în complement de doi și deci trebuie obligatoriu reconvertit în cod semn-mărime. *Așadar logica de comandă a ALU devine exagerat de complexă și conduce la întârzierea producerii rezultatului.*

Soluția pe care am propus-o și care va fi prezentată în continuare elimină aceste dezavantaje și nu utilizează decât un singur sumator pe 9 biți pentru ajustarea exponentului rezultatului.

Legătura dintre numerele care intervin în relația 2.10 și reprezentările lor în format virgulă flotantă cu exponent cumulat, marcate în figura 2.2, e dată de următoarele relații:

$$A' = A \times 2^{127}, B' = B \times 2^{127}, C' = C \times 2^{127}$$

Dacă se poate scrie $A \times B = C$, nu același lucru se întâmplă în cazul produsului $A' \times B'$. Într-adevăr $A' \times B' = A \times 2^{127} \times B \times 2^{127} = C' \times 2^{127} \neq C'$. De aici rezultă:

$$C' = 2^{(\log A' + \log B')} \times 2^{-127} \quad (2.14)$$

Așadar pentru a obține rezultatul C' plecând de la produsul $A' \times B'$ ce poate fi efectuat de către subunitatea logaritmă, este necesară scăderea din partea întregă a rezultatului furnizat de ALU a deplasamentului de 127. În felul acesta, cele două sumatoare din figura 2.2 care extrag deplasamentul din exponenții celor doi operanzi de la intrare nu mai sunt necesare.

În cazul operației de împărțire, conform relației 2.11, se poate de-asemena scrie $A / B = C$, însă se poate constata că $A' / B' \neq C'$. Într-adevăr $A' / B' = (A \times 2^{127}) / (B \times 2^{127}) = C' \times 2^{-127}$. De aici rezultă:

$$C' = 2^{(\log A' - \log B')} \times 2^{+127} \quad (2.15)$$

Cu alte cuvinte pentru a obține rezultatul C' plecând de la raportul A' / B' ce poate fi efectuat de către subunitatea logaritmă, este necesară adunarea la partea

întreagă a rezultatului furnizat de ALU a deplasamentului de 127. Nici în acest caz nu mai sunt necesare cele două sumatoare pe 8 biți.

În aceste condiții, noua arhitectură a subunității logaritmice se prezintă ca în figura 2.3.

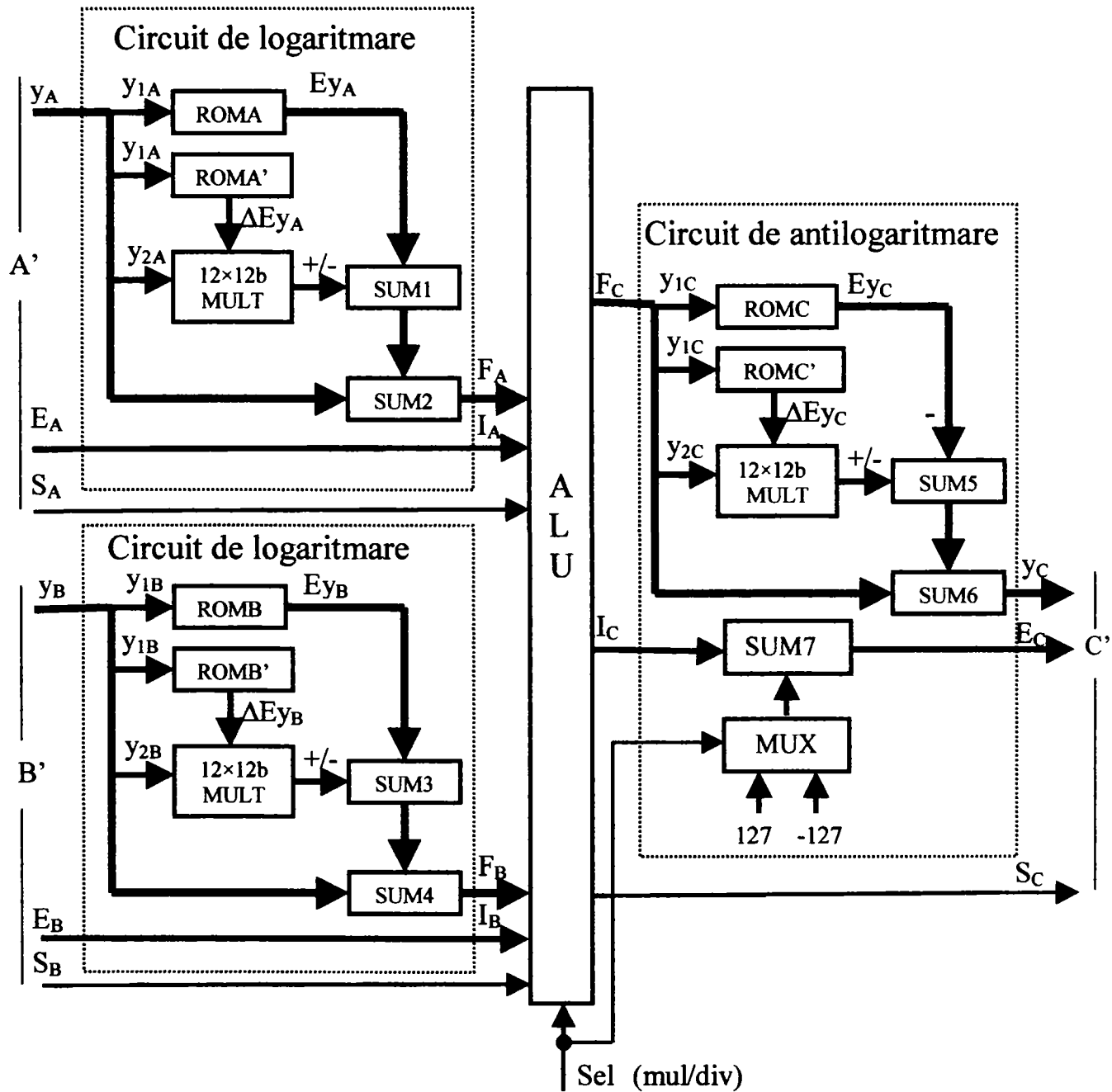


Fig.2.3. Arhitectura subunității logaritmice ce implementează algoritmul modificat

Așadar, bitul *Sel* care specifică dacă se execută multiplicarea sau divizarea operanzilor, va selecta prin intermediul blocului MUX, scăderea respectiv adunarea deplasamentului de 127 la partea întreagă a rezultatului furnizat de ALU. Acest lucru se realizează cu sumatorul SUM7.

Sumatoarele din ALU se vor extinde de la 31 de biți la 32 deoarece pot fi necesari 9 biți pentru partea întreagă a rezultatului, ținând cont că în cazul înmulțirii se adună și deplasamentele incluse în exponenții celor doi operanzi.

Acest lucru însă modifică ne semnificativ timpul de producere a rezultatului pentru sumatoarele cu transport anticipat folosite în ALU (de tip $\log(n)$, cu structură modulară arborescentă, conform prezentării din paragraful 1.1).

De asemenea SUM7 va opera pe 9 biți, dar aici întârzierea produsă la propagarea transportului datorită extinderii cu un bit, nu mai apare pe o cale critică. În plus, în cazul înmulțirii, bitul MSB (=1) de la acest sumator va fi indicator de supradepășire (când rezultă exponent cumulat mai mare de 255), iar în cazul împărțirii acest bit în conjuncție cu MSB al rezultatului furnizat de ALU va fi indicator de subdepășire (când rezultă exponent cumulat negativ), respectiv de supradepășire dacă valoarea sa logică este zero.

Eroarea maximă pe care o introduce conversia de format este 3×10^{-7} , dar în paragraful 2.2.4 autorul va arăta cum această eroare poate fi redusă cu 40%.

2.2.2 Proiectarea subunității logaritmice după metoda Lai folosind algoritmul de conversie de format modificat

În lucrarea [50] în care a fost descris algoritmul nemodificat prezentat în paragraful anterior, au fost date doar câteva informații cu privire la arhitectura unității logaritmice și la organizarea sa pe niveluri pipeline. Nu a fost făcută însă nici o specificație de proiectare cu privire la implementarea propriu-zisă a algoritmului sau la dimensionarea și programarea memoriilor, pe baza unor calcule matematice. În consecință, în acest paragraf vor fi rezolvate aceste probleme printr-o proiectare completă până la nivel de poartă logică, însă pentru varianta modificată (proprie) a algoritmului, descrisă tot în paragraful precedent. Pentru estimarea timpului total de calcul, vor fi cumulați timpii de propagare ai blocurilor din structura proiectată care se găsesc pe calea critică de propagare. Ulterior, pentru creșterea frecvenței de clock și a capacității de calcul va fi discutată posibilitatea împărțirii sale optime pe niveluri pipeline. Simulările ce vor fi efectuate pentru cazurile cele mai defavorabile din punct de vedere al propagării transportului, vor permite identificarea blocurilor ce reprezintă cale critică, în vederea repartizării cât mai uniforme a timpilor de propagare pe nivelurile pipeline.

2.2.2.1 Determinarea dimensiunii memoriilor. Proiectarea circuitului de logaritmare

Pentru a determina numărul maxim de biți diferiți de 0 pe care poate fi reprezentat E_y furnizat de memoria ROMA, în limitele formatului simplă precizie, se calculează maximul funcției $f(y) = \log_2(1+y) - y$ pe intervalul $[0,1)$. Pentru aceasta se determină rădăcina ecuației $f'(y) = 0$ pe acest interval:

$$(\log_2 e)/(1+y) - 1 = 0$$

Se obține $y=0,4427$. Pentru această valoare a lui y funcția $f(y)$ are valoarea maximă $0,086$ care reprezentată ca număr fracționar binar pe 23 de biți prezintă 3 zerouri după virgulă. Așadar, lungimea cuvintelor E_y ce trebuie memorate este $23 - 3 = 20$ biți.

Cuvintele ΔE_y reprezintă diferențele dintre valorile învecinate ale funcției $f(y)$ și se memorează la fel ca și E_y într-o tabelă de memorie de tip ROMA', care admite la intrare același cuvânt de adresă pe 11 biți. Această a doua tabelă de memorie este mult mai mică decât cea corespunzătoare lui E_y întrucât $\Delta E_y \ll E_y$ și deci numărul de zerouri după virgulă în reprezentarea fracționară binară a lui ΔE_y este mult mai mare decât 3. Pentru a determina acest număr de zerouri (care nu mai trebuie memorate) calculăm maximul funcției $f'(y)$ pe intervalul $[0,1)$. Așa cum se observă din figura 2.1.a, cel mai mare număr binar fracționar ΔE_y pe 23 de biți se obține în punctul de pe grafic unde e cea mai mare pantă, respectiv la cea mai mică adresă diferită de 0, adică 00000000001 (001h). Această valoare maximă, în sistemul de numerație binar, este:

$$\log(1+0,00000000001) - 0,00000000001 = 0,00000000000101100100111$$

și conține 11 zerouri după virgulă. Deci lungimea cuvintelor ΔE_y ce trebuie memorate în cea de-a doua tabelă de memorie ROMA' este $23 - 11 = 12$ biți.

Pentru a determina numărul maxim de biți diferiți de 0 cu care poate fi reprezentat E_y furnizat de memoria ROMC, se calculează maximul funcției $f(y)=(1+y) \cdot 2^y$ pe intervalul $[0,1)$. Pentru aceasta se determină rădăcina ecuației $f'(y)=0$ pe acest interval:

$$1 - e^y \ln 2 = 0$$

Se obține $y=0,5288$. Pentru această valoare a lui y funcția $f(y)$ are valoarea maximă $0,086$, aceeași cu cea din cazul anterior. Așadar, lungimea cuvintelor E_y ce trebuie memorate în ROMC este tot 20 biți. În mod similar se determină lungimea maximă a cuvintelor ΔE_y memorate în ROMC', observând că cea mai mare pantă pe graficul din figura 2.1.b corespunde celei mai mari adrese. Se obține o lungime tot de 12 biți.

Maniera de programare a memoriilor de tip ROMA' și ROMC' este următoarea: la adresa 0 se memorează valoarea E_y de la adresa 1 a memoriei ROMA, adică $E_y(1)$, la adresa 1 se memorează $|E_y(2)-E_y(1)|$, la adresa 2 se memorează $|E_y(3)-E_y(2)|$ ș.a.m.d. La adresa n se memorează $|E_y(n+1)-E_y(n)|$. La adresa 2047 se memorează $|0-E_y(2047)| = E_y(2047)$.

În plus, memoriile de tip ROMA' și ROMC' trebuie să memoreze în fiecare locație pe un bit suplimentar, denumit bit de control și semnul termenilor $\Delta E_y \times y_2$ din relațiile 2.6 și 2.9. Acest bit de control va avea valoarea 0 în cazul semnelui "+" și valoarea 1 în cazul semnelui "-". În cazul memoriei ROMA', în acest bit se va memora 0 de la adresa 0 la adresa 906 și 1 de la adresa 907 la adresa 2047. În

cazul memoriei ROMC' în acest bit se va memora 0 de la adresa 0 la adresa 1082 și 1 de la adresa 1083 la adresa 2047.

În concluzie, cantitatea de memorie ROM ce trebuie integrată pentru realizarea conversiilor de format, în cazul a doi operanzi de intrare și un rezultat la ieșire, este de: $2 \times [20 \times 2k + (12+1) \times 2k] + [20 \times 2k + (12+1) \times 2k] = 198$ kilobiți. Ținând cont că la fiecare din cele 3 memorii de 66kb logica de decodificare ocupă cam 60% din suprafață rezultă că în tehnologie CMOS 0,5 μm , memoriile necesare se pot integra în 1,3 mm^2 , iar timpul de acces este de 3,8 ns [50], [51].

Termenul $\Delta E_y \times y_2$ este obținut cu ajutorul unui multiplicator de 12×12 biți, astfel că numărul binar de la ieșirea acestuia va avea 24 de biți. Deoarece y_2 este văzut ca număr binar fracționar cu 12 biți după virgulă (atunci când $tg\alpha$ într-un punct pe graficul din figura 2.1.a este normal la valoarea ΔE_y) iar numărul binar fracționar ΔE_y are 23 de biți după virgulă, din care primii 11 sunt zerouri, rezultă că de la ieșirea multiplicatorului vor fi reținuți doar cei mai semnificativi 12 biți. Rezultatul propriu-zis este văzut ca un număr fracționar pe 23 biți cu 11 zerouri după virgulă, urmate de cei mai semnificativi 12 biți de la ieșirea multiplicatorului. Acest calcul a fost necesar pentru alinierea corectă a celor trei termeni care se adună din relațiile 2.6 și 2.9. Schema bloc a circuitului care realizează adunarea acestor termeni, corespunzând circuitului de logaritmare, este prezentată în figura 2.4.

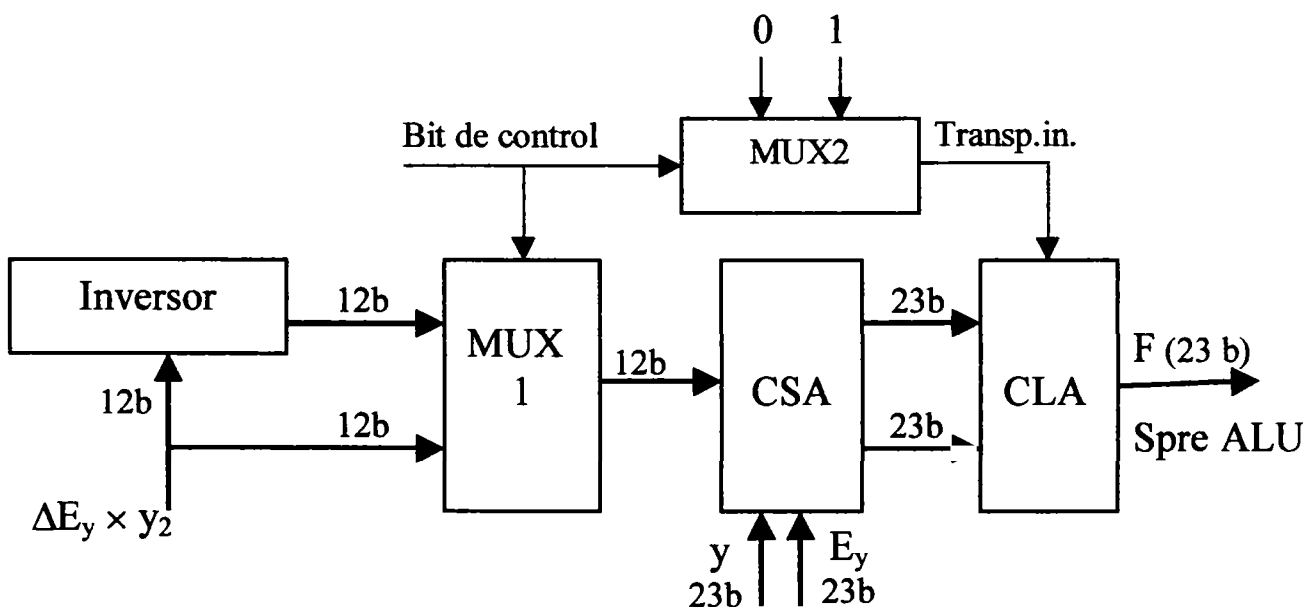


Fig. 2.4. Schema bloc a circuitului care însumează termenii y , E_y și $\Delta E_y \times y_2$

Blocul inversor asigură complementarea biților termenului $\Delta E_y \times y_2$, în vederea scăderii sale. În cazul în care ΔE_y este furnizat prin aplicarea unei adrese mai mari de 906, bitul de control are valoarea 1 și cuplează, prin intermediul blocului multiplexor MUX1, data de la ieșirea blocului inversor la

intrarea blocului CSA. În cazul în care bitul de control are valoarea 0, la intrarea blocului CSA se aplică direct termenul $\Delta E_y \times y_2$. Reprezentarea în complement de doi a termenului $\Delta E_y \times y_2$, presupune și adăugarea unui 1 biților inversați. Având în vedere comutativitatea adunării, aceasta se poate face la sumatorul final CLA (“carry look-ahead” sau cu anticiparea transportului) luând în considerare un transport inițial. Evident dacă termenul $\Delta E_y \times y_2$ trebuie adunat și nu scăzut, atunci bitul de control, prin intermediul multiplexorului MUX2, selectează transport inițial 0. Utilizarea blocului CSA (“carry save adder” sau sumator cu conservarea transportului) permite reducerea timpului de propagare pentru cele două adunări consecutive din relația 2.6 la aproape jumătate. Un segment din acest bloc este prezentat în figura 2.5:

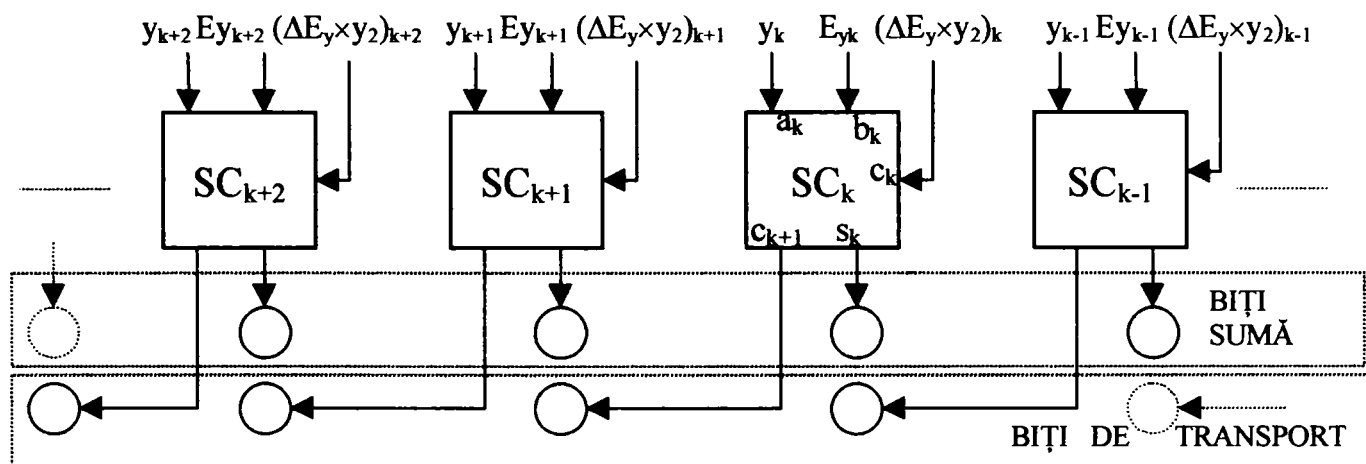


Fig. 2.5. Segment de 4 biți din blocul CSA

Blocurile SC sunt sumatoare complete pe un bit, independente, structura lor fiind arătată în figura 1.1. Biții de sumă și de transport se obțin după parcurgerea a 4 și respectiv 2 niveluri logice, adică după maxim 4 unități de timp de propagare. Astfel blocul CSA poate accepta la intrare 3 operanzi de 23 de biți și produce 2 cuvinte de câte 23 de biți, decalate cu un bit, plasate în secțiunile “sumă” și “transport”, așa cum se vede în figura 2.5.

Deoarece la nivelul unui sumator complet, din cele 23, adunarea celor 3 biți este independentă de biții de rang inferior sau superior, întreaga operație va dura 4 unități de timp de propagare. Termenii y , E_y și $\Delta E_y \times y_2$ care se aplică la intrările CSA reprezintă numere fracționare care au 23 de biți, din care al doilea are primii 3 biți după virgulă egali cu zero, iar al treilea are primii 11 biți după virgulă 0 sau 1, în funcție de valoarea bitului de control. Rezultatul final al adunării celor 3 operanzi se obține cu ajutorul unui sumator final de tip CLA care adună cuvintele din secțiunile sumă și transport ale blocului CSA. Sumatorul de tip CLA a fost descris în paragraful 1.3, dar ulterior va fi îmbunătățit.

Pentru estimarea timpului de propagare prin circuitul care însumează termenii y , E_y și $\Delta E_y \times y_2$, au fost folosite datele furnizate de [67] cu privire la timpii de propagare prin diversele porți logice corespunzătoare unei tehnologii CMOS de 0,5 micrometri. Astfel, timpul de propagare printr-o poartă SAU-

exclusiv este 0,35 ns, prin porți simple cu două intrări, de tip ȘI-NU și SAU-NU, este de 0,2 ns, în timp ce pentru un inversor sau tranzistor de trecere acest timp este de 0,1 ns. Rezultă astfel $t_p = t_{inversor} + t_{MUX} + t_{CSA} + t_{CLA} = 0,1 + 0,1 + 2 \times 0,35 + 2,3 = 3,2$ ns. Valoarea lui t_{CLA} a fost obținută prin simulare pentru cazul cel mai defavorabil al unui sumator performant de tipul celui prezentat în figura 2.12, dar pe 23 de biți.

2.2.2.2 Proiectarea multiplicatorului de 12×12 biți

Maniera asincronă în care a fost proiectat multiplicatorul de 12×12 biți a urmărit furnizarea rezultatului înmulțirii celor doi operanzi într-un timp cât mai scurt. Din acest motiv au fost folosite blocuri CSA, similare cu cel descris în paragraful anterior, care grupează produsele parțiale din aria de înmulțire într-un arbore binar de tip Wallace. Viteza mare de calcul se obține pe de o parte datorită propagării doar pe verticală a transportului prin blocurile CSA din aria de înmulțire, iar pe de altă parte datorită calculării simultane a mai multor produse parțiale intermediare, avantaj oferit de structura în arbore Wallace. Schema bloc a acestui arbore este prezentată în figura 2.6.

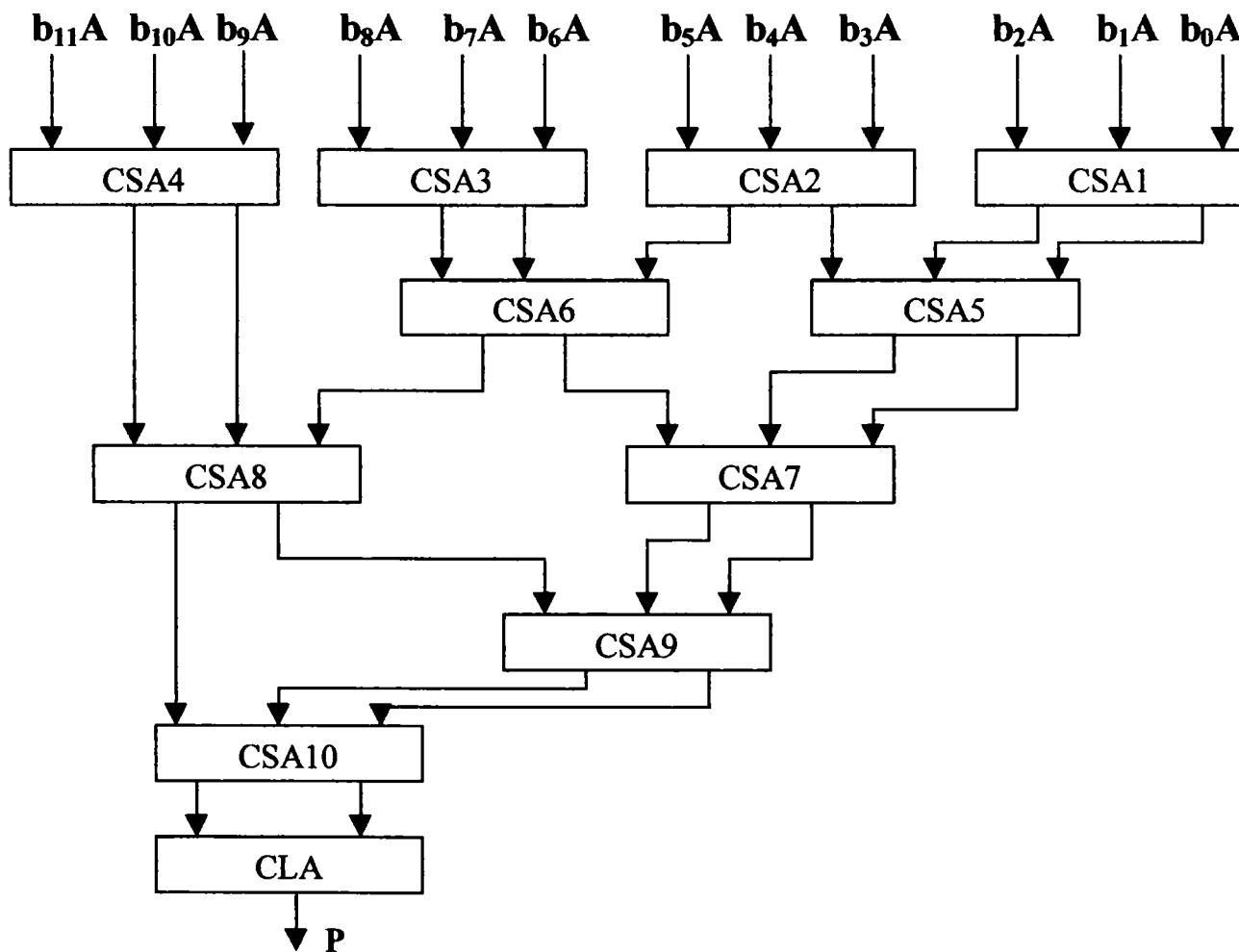


Fig. 2.6. Schema bloc a arborelui Wallace realizat cu blocuri CSA

Așa cum se observă în această figură, fiecare bloc CSA acceptă la intrare trei produse parțiale inițiale sau trei produse parțiale intermediare, furnizând la ieșire alte două produse parțiale intermediare. Are loc astfel o reducere treptată a

numărului de produse parțiale, de-a lungul a 5 niveluri de blocuri CSA, până când, în final, rezultă doar două cuvinte ce trebuie însumate. Toate blocurile CSA de pe același nivel lucrează simultan. Fiind vorba de structuri cu conservarea și nu propagarea pe orizontală a transportului, un nivel CSA este parcurs în maxim 4 unități de timp de propagare. Aceasta înseamnă că timpul de propagare prin cele 5 niveluri este de 20 unități de timp de propagare.

Deoarece produsele parțiale inițiale, au fiecare pondere dublă față de cel anterior (trebuie aliniate decalat cu câte un bit), iar diferitele produse parțiale intermediare care se grupează într-un bloc CSA au ponderi diferite (se aliniază decalat, în funcție de ponderea fiecăruia), rezultă lungimi diferite pentru blocurile CSA utilizate. De exemplu, în blocul CSA5 ponderile și lungimile celor trei produse parțiale intermediare care se grupează aici au condus la alinierea prezentată în figura 2.7. Lungimea acestui bloc este de 15 biți, din care bitul LSB din secțiunea biților sumă de la ieșirea sa, ne mai combinându-se ulterior cu nimic, ajunge direct bitul P_2 al produsului final. Se observă de asemenea că în acest caz nu sunt necesare decât 11 sumatoare complete pe 1 bit și un sumator incomplet, întrucât cei mai semnificativi trei biți din al treilea produs parțial intermediar de intrare ajung direct în secțiunea de ieșire a biților sumă, de lungime 14 biți. Secțiunea biților de transport conține aici 12 biți.

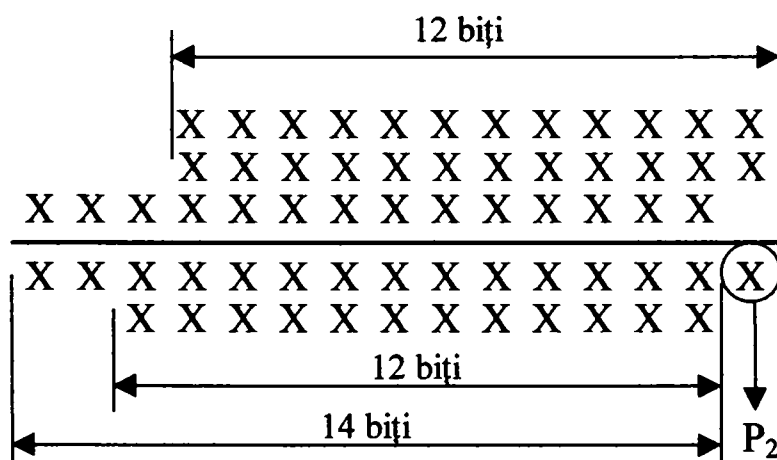


Fig. 2.7. Alinierea produselor parțiale intermediare în blocul CSA5

Spre vârful arborelui produsele parțiale intermediare devin din ce în ce mai lungi, dar aceasta nu conduce la creșterea semnificativă a numărului de porți logice utilizate într-un bloc CSA deoarece un grad mai mic de suprapunere a cuvintelor de intrare presupune utilizarea mai multor sumatoare incomplete în loc de sumatoare complete pe un bit, sau poate însemna chiar transferul biților mai semnificativi direct la ieșire.

Așa cum s-a putut constata și în cazul blocului CSA5, cei mai puțin semnificativi 6 biți ai produsului sunt furnizați de structura hard a arborelui Wallace, astfel că lungimea necesară a sumatorului final este de 18 biți. S-a optat pentru un sumator rapid cu anticiparea transportului (de tip CLA). Acesta a fost proiectat conform regulilor descrise în paragrafele 1.1 și 2.2.2.4. Pentru calculul produselor parțiale inițiale de forma $b_i A$, cu $i = 0 \div 11$ au fost utilizate blocuri de

câte 12 porți ȘI cu 2 intrări și arbori de inversoare cu 2 etaje (cu 2 și respectiv 4 inversoare pe nivel) pentru multiplicarea capacității de comandă a biților celor doi factori, în vederea generării celor 12 produse parțiale. Timpul de propagare prin multiplicator rezultă: $t_p=2 \times t_{inv} + t_{\text{ȘI}} + t_{\text{arbore}} + t_{\text{CLA}} = 0,2 + 0,3 + 10 \times 0,35 + 2,3 = 6,3 \text{ ns}$.

2.2.2.3 Circuitul de antilogaritmare

Proiectarea secțiunii circuitului de antilogaritmare care se ocupă de partea fracționară a rezultatului furnizat de ALU are ca bază de plecare relația 2.9. Diferențele, din punct de vedere formal, față de relația 2.6 care a condus la implementarea circuitului de logaritmare, constau în acelea că termenul E_y trebuie scăzut și nu adunat și că apare valoarea 1 în paranteza $(1+y)$. Această valoare este memorată implicit în structurile hard care operează cu numere în virgulă flotantă normalizate, deci nu e necesar ca ea să se regăsească pe undeva în structura circuitului de antilogaritmare, acesta furnizând implicit o mantisă normalizată.

Se observă din relația 2.9 că dacă se înlocuiește scăderea termenului E_y cu însumarea complementului de doi al acestuia, aceasta conduce la o implementare hard identică a circuitului de antilogaritmare cu cea a circuitului de logaritmare prezentat la nivel de schemă bloc în figura 2.4.

Diferențele se regăsesc evident în conținutul memoriilor ROMC și ROMC'. Astfel ROMC va memora complementul de doi al cantităților E_y corespunzătoare reprezentării grafice din figura 2.1.b în loc de chiar E_y . În cazul memoriei ROMC', programarea sa se va face similar cu a memoriei ROMA' pe baza valorilor E_y precalculate, corespunzătoare graficului din figura 2.1.b. În locațiile corespunzătoare bitului de control al semnului cantității $\Delta E_y \times y_2$ se va memora 0 de la adresa 0 la adresa 1082 și 1 de la adresa 1083 la adresa 2047.

În ceea ce privește secțiunea circuitului de antilogaritmare, care furnizează exponentul cumulat al rezultatului, aceasta se regăsește în figura 2.3. Sumatorul SUM7 din figura 2.3 care produce pe E_C este prezentat în figura 2.8.

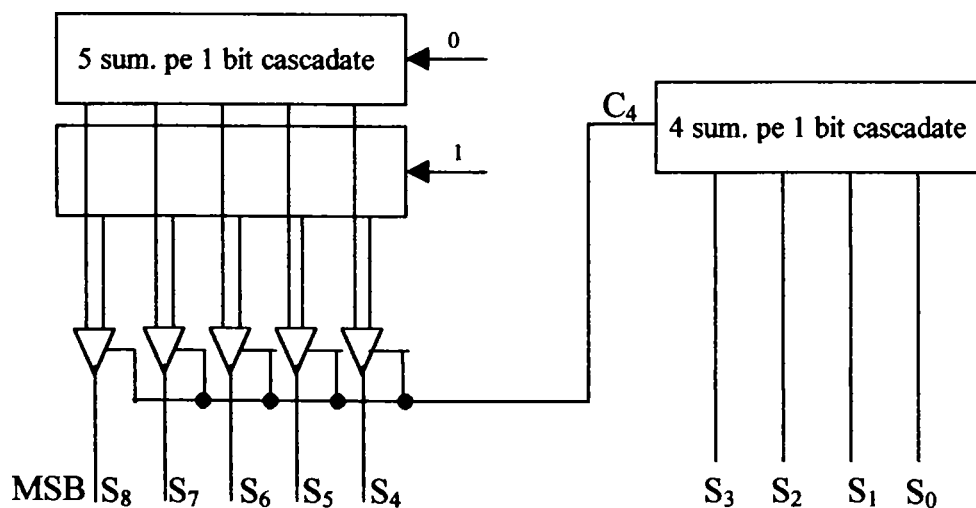


Fig.2.8. Schema bloc a sumatorului care produce exponentul rezultatului

SUM7 este un sumator cu selecția transportului organizat pe două secțiuni de 4 și respectiv 5 biți. Transportul C_4 selectează rezultatul corect antecalculat de cele două sumatoare pe 5 biți ce primesc transport inițial 0, respectiv 1. Cele două secțiuni includ sumatoare complete pe un bit (sumatoare complete 3:2), cascade, la care întârzierea produsă de propagarea din aproape în aproape a transportului se face pe o cale care nu este critică, întrucât SUM7 funcționează în paralel cu cealaltă parte a circuitului de antilogaritmare, care furnizează pe y_C într-un interval de timp mai mare. Ieșirea circuitului de antilogaritmare reprezintă practic ieșirea subunității logaritmice.

2.2.2.4 Proiectarea circuitului de adunare/scădere a operanzilor în format LNS, cu deplasament

Cei doi operanzi în format LNS simplă precizie și cu deplasament, furnizați de cele două circuite de logaritmare ce funcționează simultan, au așadar lungimea de 32 de biți. Reprezentarea lor în virgulă fixă include 9 biți pentru partea întreagă și 23 de biți pentru partea fracționară. Din punct de vedere hard, circuitul care realizează adunarea sau scăderea lor este identic cu acela care realizează adunarea sau scăderea a doi operanzi întregi de 32 biți.

Includerea deplasamentului în partea întreagă a operanzilor de la intrarea în ALU face ca aceștia să fie întodeauna pozitivi. În aceste condiții structura ALU se simplifică, ea se reduce la un circuit de adunare/scădere care operează cu date pozitive.

Efectuarea operației de împărțire în subunitatea logaritmică presupune scădere în ALU, adică conversia în complement de doi a scăzătorului. Dacă, în plus, rezultatul se dovedește a fi negativ, atunci el este automat reprezentat în complement de doi și trebuie ulterior reconvertit în cod semn-mărime. Întârzierea produsă ar putea fi eliminată, dacă pentru cei doi operanzi A și B de la intrarea în ALU se calculează simultan $A-B$ și $B-A$ după care se selectează rezultatul pozitiv și se stabilește semnul corect. Și în acest caz trebuie asigurată conversia în complement de doi pentru ambii operanzi cât și posibilitatea efectuării adunării propriu-zise. Soluția pe care am adoptat-o, [26], se bazează pe următoarele relații:

$$A - B = A + \overline{B} + 1 = A + /B + 1 \quad (2.16)$$

$$B - A = \overline{\overline{A + B}} \quad (A - B = -\overline{\overline{B - A}} = \overline{\overline{A + B + 1}}) \quad (2.17)$$

în care $/B$ este data obținută prin inversarea biților lui B. Pentru efectuarea operației de scădere, circuitul sumator/scăzător calculează simultan $A+/B$ și $A+/B+1$ utilizând două sumatoare pe 32 de biți, unul fără transport inițial iar celălalt cu transport inițial. Dacă $A+/B+1$ este pozitiv, el este selectat la ieșirea ALU. Dacă $A+/B+1$ este negativ, atunci $A+/B$ este selectat la ieșire, după ce în prealabil biții săi au fost inversați. Schema bloc a acestui circuit este prezentată în figura 2.9. Spre deosebire de schema propusă în [26], soluția proprie are

topologia blocurilor componente schimbată, permițând selectarea corectă a rezultatului, nu numai în cazul scăderii, dar și în cazul adunării operanzilor. În plus, schema propusă e mai rapidă deoarece blocul Inversor2 nu mai așteaptă stabilirea bitului MSB, marcat în figură și apoi rezultatul de la blocul Selector. În funcție de valoarea bitului MSB, dar condiționat de linia de control Sel (Sel=0 în cazul adunării, respectiv Sel=1 în cazul scăderii) este selectat rezultatul de la blocul Sumator1, inversat sau neinversat, sau de la blocul Sumator2.

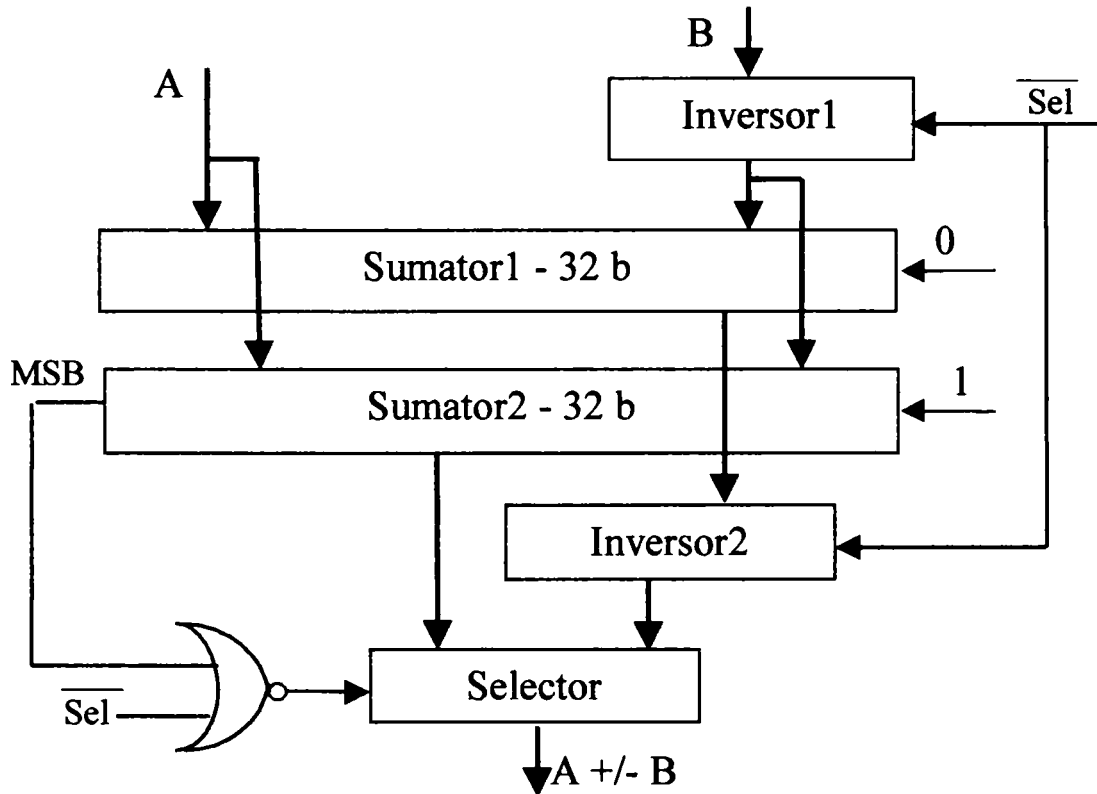


Fig. 2.9. Schema bloc a circuitului de adunare/scădere din componența ALU

Atunci când se efectuează adunare, doar blocul Sumator1, cu transport inițial 0, operează, în timp ce blocurile Inversor1 și Inversor2 vor fi transparente.

Cele două sumatoare sunt sumatoare rapide, alese într-o primă etapă, de tip CLA (carry look-ahead) sau cu anticiparea transportului și au fost prezentate în paragraful 1.1. Dacă se analizează schema la nivel de poartă a blocurilor de tip A și B prezentate în figurile 1.5 și 1.6 care formează structura modulară a sumatorului, se constată că sunt utilizate porți ȘI, porți SAU și porți SAU-exclusiv. Dacă în tehnologie CMOS poarta SAU-exclusiv, așa cum se va vedea într-un alt paragraf, are o structură optimizată, nu același lucru se întâmplă cu porțile ȘI și SAU care practic se realizează prin adăugarea a câte unui inversor CMOS la porțile ȘI-NU, respectiv SAU-NU intrinseci. Aceasta ar însemna, pentru un sumator rapid pe 32 de biți, integrarea unui număr de 219 inversoare, care conduce evident la creșterea suprafeței ocupate, cât și a timpului de operare.

Întrucât literatura de specialitate nu prezintă soluții concrete pentru implementarea optimă în tehnologie CMOS a unui sumator rapid de tipul carry look-ahead (sau cu anticiparea transportului) cu structură modulară, de tipul celei

din figura 1.4, am oferit o soluție personală pentru eliminarea dezavantajelor menționate.

Astfel, am modificat structura blocurilor de tip A și B care se regăsesc în schema bloc a sumatorului din figura 1.4. Noul bloc A' derivat din blocul de tip A va avea structura din figura 2.10.a. Acest bloc va furniza practic coeficienții p și g din relațiile (1.3) complementați. Introducerea, într-o etapă intermediară de sinteză, a unor inversoare pe aceleași linii, pentru refacerea logicii, dar în interiorul blocurilor de tip B aflate imediat sub blocurile A' și aplicarea apoi a transformărilor De Morgan, a condus la obținerea blocurilor de tip B' a căror structură este prezentată în figura 2.10.b. Deoarece starea logică de la ieșirile blocurilor B' este identică cu cea de la ieșirea blocurilor de tip B, a doua linie de blocuri de tip B din schema inițială, se transformă ca și blocurile A, pentru a furniza ieșiri complementate, obținându-se blocurile de tip B'' (figura 2.10.c).

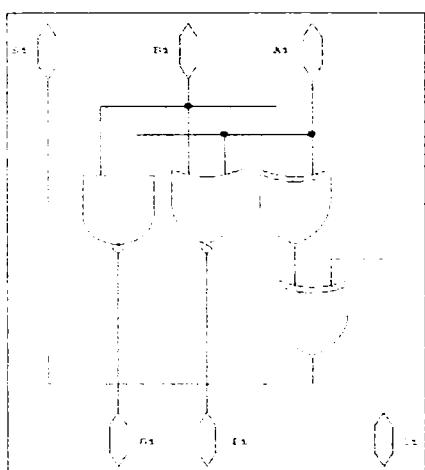


Fig.2.10.a
Schema blocului de tip A'

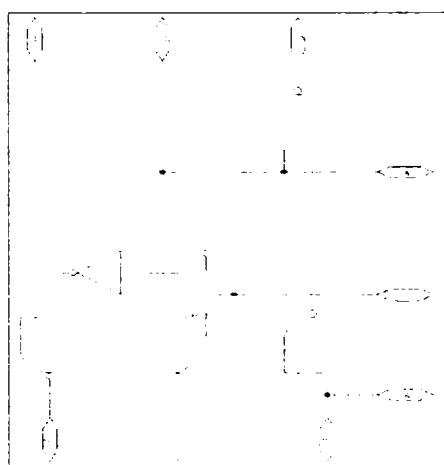


Fig.2.10.b
Schema blocului de tip B'

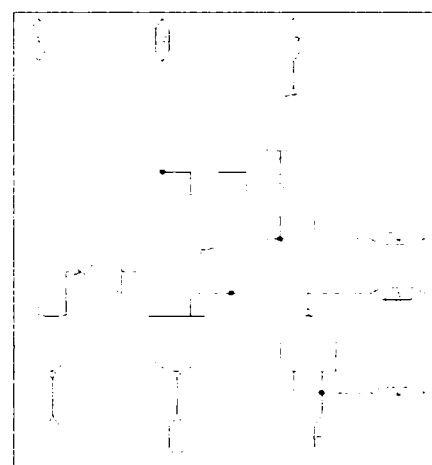


Fig.2.10.c
Schema blocului de tip B''

În acest fel se asigură compatibilitatea logică cu a treia linie de blocuri care va fi din nou de tip B'. A patra linie de blocuri va fi de tip B'', păstrându-se astfel alternanța. Schema bloc a sumatorului pe 32 de biți menține așadar aceeași organizare, dar conformația blocurilor componente și într-o anumită măsură și disponerea lor va fi diferită, așa cum se observă în figura 2.11.

Schema internă la nivel de poartă logică a acestui sumator nu mai conține decât porți optimale pentru tehnologie CMOS, respectiv porți cu două intrări de tip ȘI-NU, SAU-NU și SAU-exclusiv, respectiv inversoare CMOS. Numărul de inversoare necesar a fi integrate s-a redus astfel de la 219 la 62, iar câștigul de viteză obținut, așa cum se va vedea în capitolul ce prezintă rezultatele simulărilor, a fost de peste 25%. Această procedură va fi extinsă la orice alt sumator de același tip din unitatea de calcul. Au fost făcute simulări cu programul Msim52, la care timpii de propagare ai diverselor porți logice utilizate, extrase din biblioteca de modele a simulatorului, au fost setați în conformitate cu datele oferite de [67], corespunzătoare unei tehnologii CMOS de 0,5 μ m, adoptată ca și etalon.

Conceptia subunității logaritmice a unui procesor aritmetic hibrid

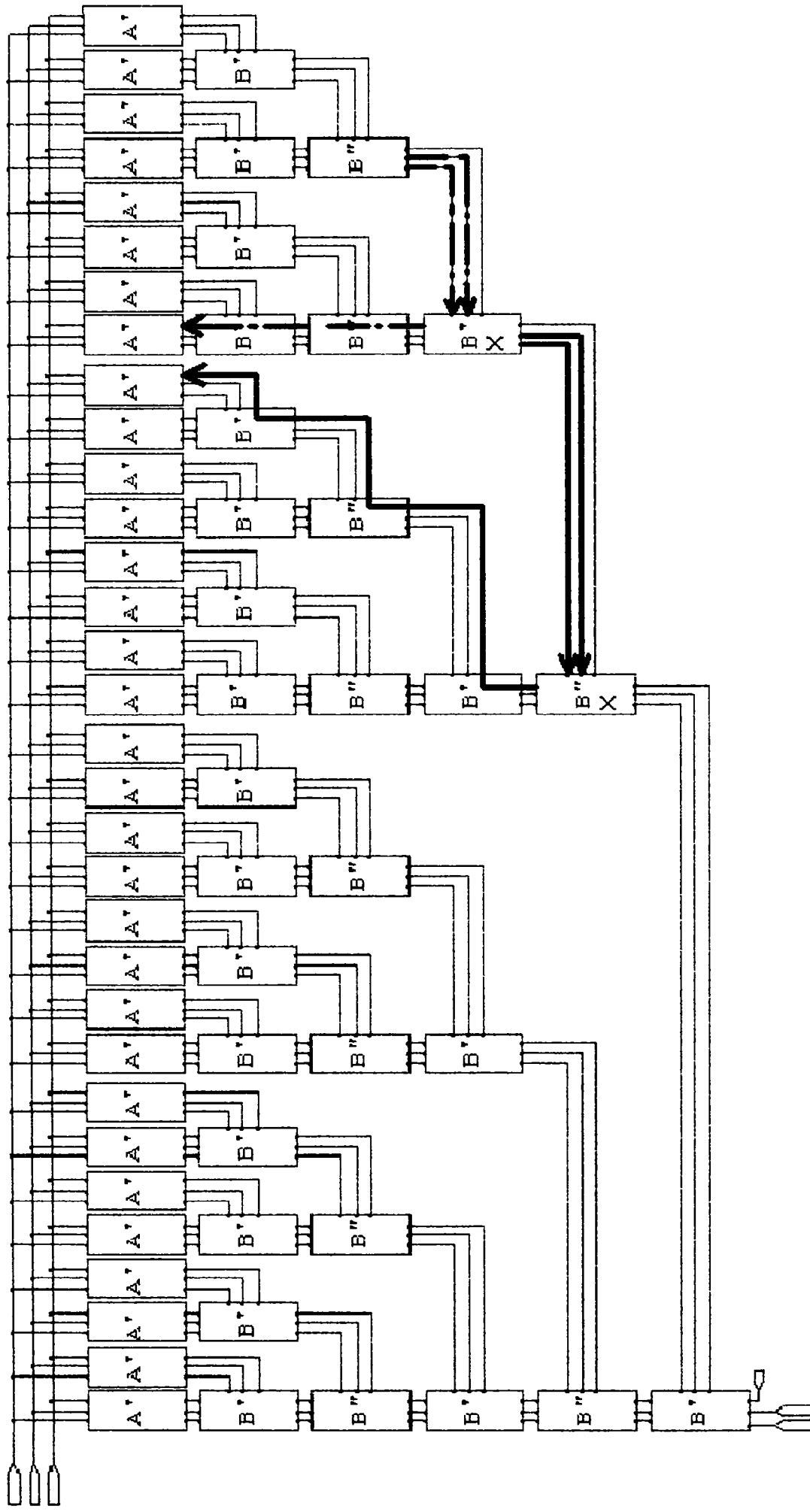


Fig. 2.11. Schema bloc a sumatorului CLA cu structură modulară, pe 32 de biți, modificat

Pentru sumatorul de tip CLA pe 32 de biți, în urma modificărilor efectuate, timpul de calcul a coborât de la 5,7 ns la 4,3 ns, pentru cazul cel mai defavorabil. În cazul unui sumator de tip CLA pe 16 biți, timpul de calcul a fost de 3,5 ns, iar pentru unul pe 8 biți, acest timp a fost de 2,4 ns.

Depistarea cazului celui mai defavorabil, respectiv a perechii de operanzi, pentru care se obține cel mai mare timp de propagare, la un sumator CLA pe 8 biți, am realizat-o cu ajutorul unui program scris în limbaj C. Programul implementează relațiile recurente (1.5), (1.6) și (1.7) și asociază contoare care se incrementează la fiecare schimbare a stării logice după fiecare nivel logic parcurs. El este prezentat în Anexa2. În urma rulării programului au fost testate toate cele 65536 combinații posibile de operanzi și au fost găsite 256 de perechi de valori pentru care se obține timpul maxim de generare a rezultatului. Între aceste perechi de valori s-a găsit și cea care corespunde cazului celui mai defavorabil de la sumatorul cu propagarea din aproape în aproape a transportului (la care 8 sumatoare complete pe un bit sunt cascadeate), respectiv FFh și 01h. În schimb cazul FFh și 00h cu transport inițial 1 nu reprezintă un cel mai defavorabil caz. Acest lucru a fost confirmat și prin simulare. Concluzia importantă care s-a desprins a fost că, datorită structurii modulare a unui sumator CLA, cazul particular obținut a putut fi extins ca și cel mai defavorabil caz și la sumatoare CLA mai mari. Pentru un sumator CLA pe 32 de biți aceasta a însemnat FFFFFFFFh și 00000001h.

Etapa următoare în proiectare, pe care am abordat-o, a fost analiza generării și propagării informației de la primul bloc de 8 biți al sumatorului CLA pe 32 de biți, în vederea stabilirii transportului inițial al următorului bloc de 8 biți. Observația importantă făcută a fost că transportul inițial (marcat c8 în figura 2.11) pentru al doilea bloc de 8 biți se obține cu un interval de timp semnificativ înaintea unora din primii 8 biți ai sumei. Într-adevăr, așa cum se observă din figura 2.11, calea pentru generarea lui c8, de la ieșirea blocului B'x, marcată cu linie continuă îngroșată, nu presupune decât parcurgerea a două niveluri logice în blocul B''x. Celelalte blocuri de pe traseu vor fi transparente pe linia menționată, așa cum se observă din figurile 2.10 și 2.11. În schimb, calea pentru generarea lui c7, necesar pentru calcularea bitului s7 al sumei, marcată cu linie întreruptă îngroșată, presupune parcurgerea a 4 niveluri logice prin blocurile B'' și B' de deasupra lui B'x. La acestea se adaugă încă două niveluri logice, în blocul A' pentru generarea lui s7.

Această proprietate mi-a permis să concep un tip nou de sumator rapid pe 32 de biți care în ansamblu este un sumator cu selecția transportului, de tipul celui prezentat în figura 1.8, dar la care sumatoarele componente, pe câte 8 biți, sunt de tipul CLA. Schema bloc a acestui sumator este prezentată în figura 2.12.

Blocurile de tip CONECTOR1,2 permit doar disocierea celor două magistrale de intrare de 32 de biți în câte 4 magistrale de 8 biți. Blocul CONECTOR3 asociază cele 4 magistrale de ieșire de 8 biți într-una singură de 32 de biți. Circuitul conține un prim sumator CLA cu transport inițial 0 și trei perechi

de sumatoare CLA cu transport inițial 0 (blocurile "SUM8_Cin=0"), respectiv 1 (blocurile "SUM8_Cin=1"). Toate sumatoarele lucrează simultan. În momentul în care transportul c8 de la ieșirea primului sumator este cunoscut, el va selecta prin intermediul blocului SELECTOR1 rezultatul de la acel sumator din prima pereche care a avut transportul inițial egal cu c8. În același timp este selectat și transportul de ieșire corect, furnizat de același sumator. Acest transport va selecta apoi, prin intermediul blocului SELECTOR2, rezultatul corect împreună cu transportul de ieșire asociat, din a doua pereche de sumatoare. Deoarece transporturile c8 se obțin înaintea stabilirii tuturor biților de ieșire ai sumatoarelor, cele 4 grupe de câte 8 biți ai sumei finale se obțin aproape instantaneu.

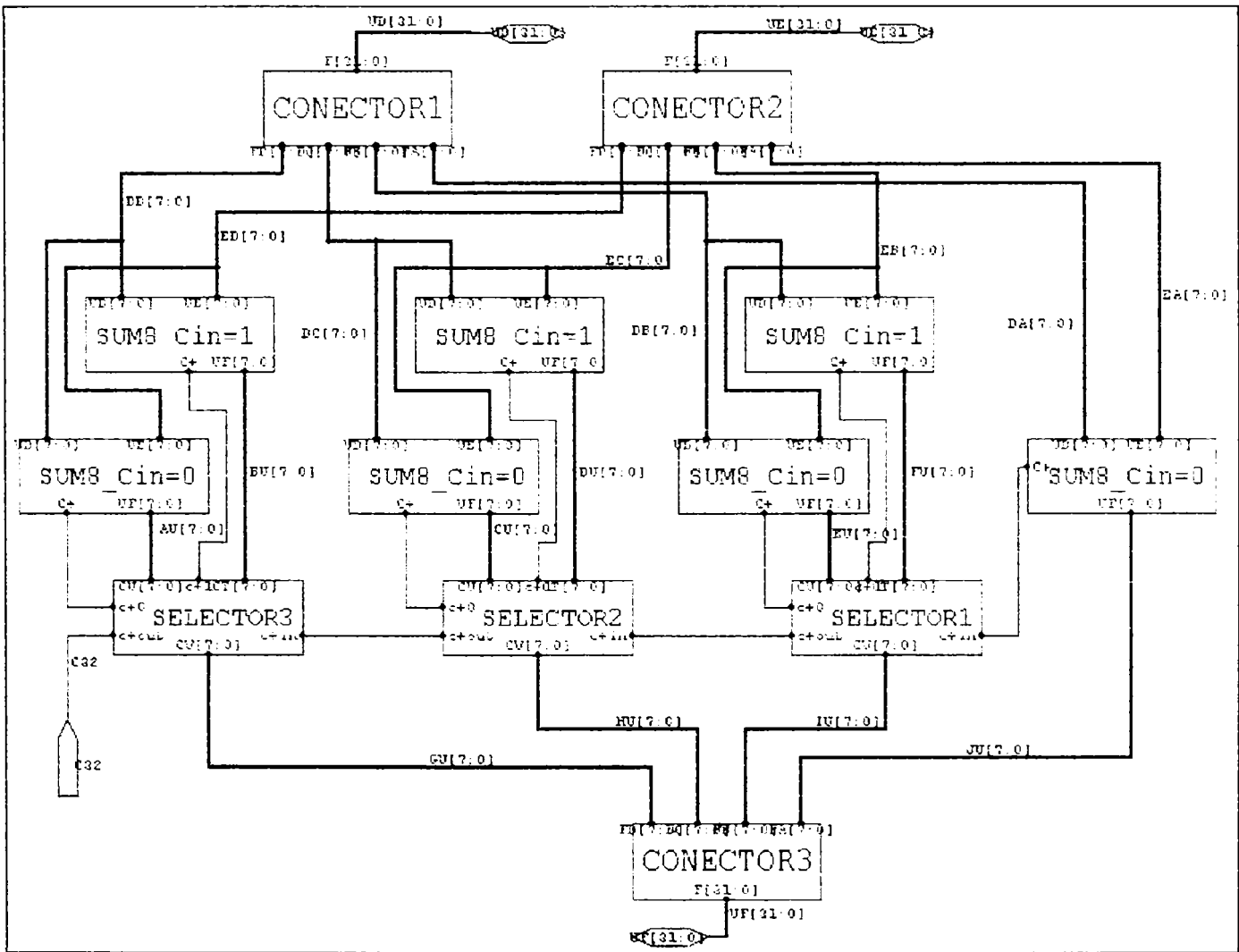


Fig. 2.12. Schema bloc a noului sumator rapid pe 32 de biți

Blocurile selectoare sunt realizate cu câte 9 perechi de tranzistoare de trecere complementare. În realitate, fiecare bit c8 trebuie să comande în fiecare bloc selector 18 grile de tranzistoare MOS. Multiplicarea capacității liniilor de selecție cu ajutorul unor arbori cu inversoare CMOS ar conduce însă la o întârziere cumulată semnificativă. Acest neajuns este compensat în bună parte datorită proprietății evidențiate anterior a sumatorului de tip CLA. În plus, am definit o linie de selecție prioritară pentru selectarea transportului inițial corect

pentru următoarea pereche de sumatoare, arborele de inversoare dezvoltându-se în aval de aceasta. Acest lucru poate fi observat în figura 2.13. Aici, transportul inițial "C+in" selectează "C+out" dintre "C+0" și "C+1" deja produși și ulterior se multiplică capabilitatea liniei de selecție. Astfel se menține practic viteza sumatorului pe 32 de biți ca și a unuia pe 8 biți.

Prin mecanismul de selecție adoptat, fiecărei ieșiri CMOS îi revine un număr de maxim 4 grile MOS pentru comandă. Deoarece arborele de inversoare are un număr impar de etaje, logica de selecție a blocului multiplexor MuxTrib este inversă față de cea a blocurilor Mux2-1.

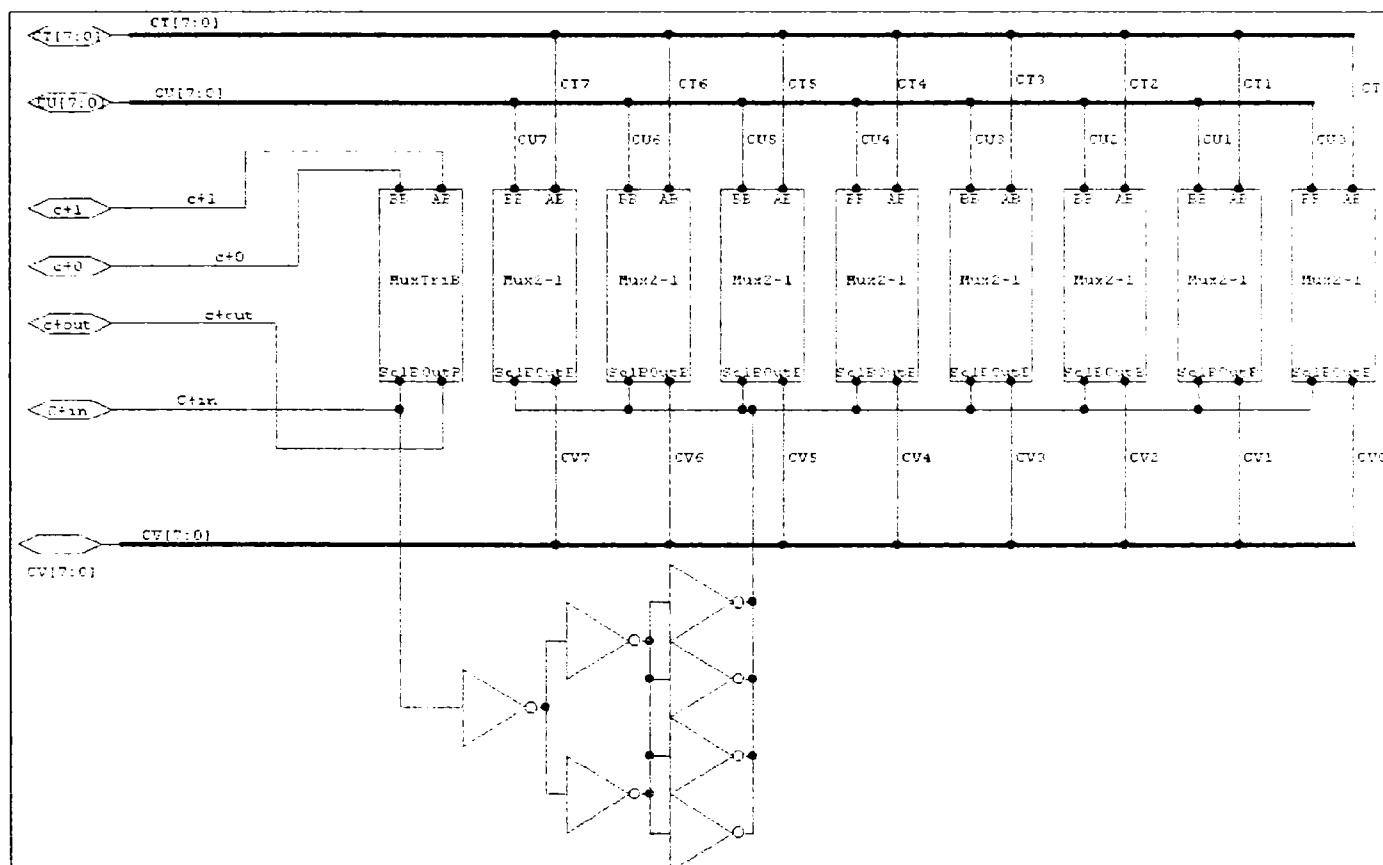


Fig. 2.13. Schema blocului Selector ce include calea de selecție prioritară

În urma simulărilor, pentru acest sumator hibrid s-a confirmat un timp de propagare de 2,6 ns, pentru cazul cel mai defavorabil. Viteza obținută este deci cu aproape 70% mai mare față de a sumatorului CLA modificat din figura 2.11 și cu 120% mai mare decât a sumatorului CLA nemodificat de tipul celui din figura 1.4, extins la 32 de biți. Aria ocupată este însă mai mare cu aproximativ 75% față de a sumatorului CLA modificat.

Făcând o comparație de viteză cu unele realizări recente în domeniu, sintetizată prin rezultatele din tabelul 2.1, se poate constata viteza superioară a sumatorului propriu. Toate valorile sunt raportate la tehnologia considerată etalon în această lucrare și anume 0,5 micrometri CMOS.

Plecând de la acest tip de sumator a fost proiectat circuitul sumator/scăzător pe 32 de biți a cărui schemă bloc a fost prezentată în figura 2.9. Schema sa detaliată, organizată pe 4 niveluri ierarhice, este prezentată în Anexa 1.

Pe lângă acest circuit, blocul ALU din figura 2.3 mai include poarta SAU-exclusiv pentru stabilirea semnului rezultatului înmulțirii/împărțirii celor doi operanzi de la intrarea subunității logaritmice de calcul.

Tab. 2.1

Sumator 32 biți	propriu	[50] ₁₉₉₁	[85] ₁₉₉₆	[78] ₁₉₉₉
Timp de propagare	2,6 ns	5ns	4,2ns	3ns

Simulările efectuate pentru circuitul sumator/scăzător, pentru cazul cel mai defavorabil, au confirmat viteza foarte mare de calcul. Timpul maxim de propagare a fost găsit de 3.3 ns, în tehnologie CMOS 0,5μm. Deoarece însă la blocul Inversor1 din figura 2.9 linia "/Sel" este deja setată (odată cu stabilirea tipului operației pe care o execută subunitatea logaritmică) înainte ca circuitele de logaritmare din figura 2.3 să furnizeze cei doi operanzi pentru ALU, întârzierea produsă la multiplicarea capabilității de comandă a liniei /Sel, de 0,5 ns, este eliminată. În concluzie, timpul maxim de operare pentru ALU este de 2,8 ns. Și în acest caz se poate constata viteza superioară circuitului propriu, față de unele realizări recente reprezentative marcate în tabelul 2.2.

Tab.2.2

Ad/Scăd. 32 biți	propriu	[26] ₁₉₉₂	[19] ₁₉₉₉	[32] ₂₀₀₀
Timp de propagare	2,8 ns	10ns	4,6ns	3,3ns

Viteza superioară a noului circuit de adunare/scădere este datorată, pe de o parte utilizării noului tip de sumator pe 32 de biți, iar pe de alta, topologiei îmbunătățite a circuitului de adunare/scădere.

2.2.2.5 Concluzii la paragraful 2.2.2

Calculul timpului total de producere a rezultatului se va face prin cumularea timpilor de propagare ai tuturor blocurilor aflate pe calea critică de propagare. Se obține astfel:

$$\begin{aligned} t_{\text{tot}} &= t_{\text{ROMA}} + t_{\text{MUL.log}} + t_{\text{CSA+CLA.log}} + t_{\text{ALU}} + t_{\text{ROMC}} + t_{\text{MUL.alog}} + t_{\text{CSA+CLA.alog}} = \\ &= 3,8 \text{ ns} + 6,3 \text{ ns} + 3,2 \text{ ns} + 2,8 \text{ ns} + 3,8 \text{ ns} + 6,3 \text{ ns} + 3,2 \text{ ns} = 29,4 \text{ ns} \quad (2.18) \end{aligned}$$

Dacă întreaga structură din figura 2.3 ar funcționa într-o manieră complet asincronă, deci dacă ar fi organizată într-un singur etaj, atunci frecvența de clock a sistemului nu va putea fi mai mare de 33 MHz. Din acest motiv s-a adoptat soluția practică în toate procesoarele aritmetice sau de semnal actuale, aceea de organizare a întregii structuri pe niveluri pipeline, așa cum se vede în figura 2.14.

Structura prezentată în figura 2.14 reprezintă varianta optimă din punct de vedere al segmentării ei în scopul asigurării frecvenței de clock mărite și conține 6 niveluri pipeline.

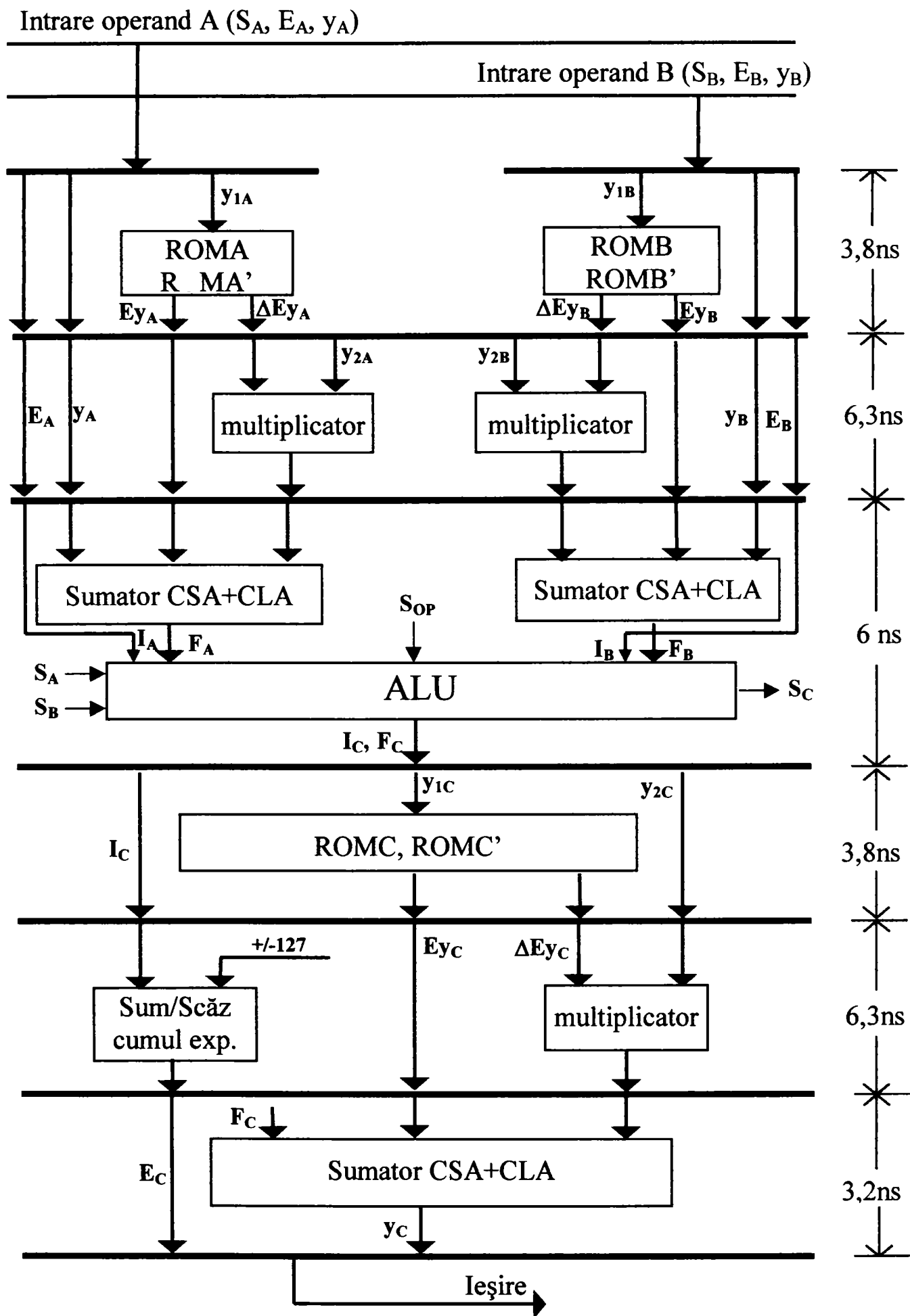


Fig. 2.14 Organizarea subunității logaritmice pe niveluri pipeline

Etajul critic al structurii este cel care conține multiplicatorul (6,3 ns), deci sistemul va putea funcționa cu o frecvență de clock maximă de aproximativ 150 MHz. Blocurile “Sumator CSA+CLA” și ALU au fost comasate în același etaj, întrucât timpul lor cumulativ de propagare este mai mic decât cel al multiplicatorului de 12×12 biți. Rezultatul unei înmulțiri sau împărțiri va fi furnizat cu o latență de 6 cicluri de 6,3 ns ($t_L = 38$ ns), dar rezultatul unor înmulțiri sau împărțiri succesive de operanzi se va produce la fiecare clock.

Rezultatele obținute în ce privește viteza de calcul sunt similare cu cele prezentate în [50] dar numărul de porți logice utilizate este mai mic deoarece au fost eliminate două sumatoare de câte 8 biți (pentru extragerea deplasamentului din exponenții celor doi operanzi), precum și două sumatoare pe 23 de biți ce calculează termenii $y+E_y$ pentru cei doi operanzi și s-a simplificat logica de comandă a ALU.

Comparativ cu [67] la care o înmulțire în virgulă flotantă se face în 10 ns, în aceeași tehnologie, se poate spune că o operație singulară se efectuează de patru ori mai lent, dar în contextul utilizării procesorului în aplicații DSP sau grafice, transformări de coordonate etc., atunci când se efectuează un număr foarte mare de produse consecutive (când un produs se obține la fiecare 6,3 ns), performanțele obținute sunt superioare sistemelor de calcul în virgulă flotantă, dotate cu un multiplicator performant de tipul prezentat în [67]. *Avantajul provine din posibilitatea mult mai facilă de segmentare a subunității logaritmice pe niveluri pipeline și deci de asigurare a unei frecvențe de clock mai mari. În plus operația de împărțire, chiar și efectuată singular, cu o latență de șase cicluri, este mai rapidă decât una efectuată în virgulă flotantă [53], [54] cu cel puțin 50%.*

2.2.3 Implementarea algoritmului modificat de conversie de format FP-LNS și LNS-FP prin metoda logaritmului mascat

2.2.3.1 Metoda logaritmului mascat

Structura prezentată în figura 2.14, deși optimizată din punct de vedere al împărțirii pe niveluri pipeline, prezintă totuși câteva dezavantaje:

- datorită funcționării sincrone pe cele șase niveluri, se adaugă aproape 9 ns la timpul total de producere a rezultatului în manieră de funcționare pur asincronă. Într-adevăr $t_L - t_{tot} = 38 - 29,4 = 9,2$ ns;

- avantajul de viteză adus prin propagarea pe verticală a transportului în aria de înmulțire organizată în arbore Wallace, este diminuat de timpul necesar pentru însumarea propriu-zisă a ultimelor două produse parțiale intermediare (când transportul se propagă pe orizontală), chiar dacă aceasta se face cu un sumator rapid;

- în nivelul al treilea al structurii, există două etape succesive de însumare propriu-zisă ce concură la obținerea unui rezultat final unic, la care se pierde mult timp prin propagarea pe orizontală a transportului;

- avantajul de viteză adus prin concepția sumatorului propriu, prezentat în paragraful 2.2.2.4, nu poate fi pe deplin valorificat datorită faptului că nivelul critic al structurii pipeline rămâne cel care conține multiplicatorul;

- multiplicatorul de 12×12 biți a fost proiectat fără a utiliza și blocuri de compresori 4:2 pe lângă blocuri CSA de sumatoare complete 3:2, ca și metoda eficientă de reducere a numărului de produse parțiale (dar cu consum suplimentar de arie pe cip, deoarece un compresor 4:2 conține două sumatoare complete 3:2); altfel, pentru structura din figura 2.14 nivelul al treilea devenea critic și câștigul de viteză ar fi fost puțin semnificativ în raport cu aria suplimentară necesară.

Pentru eliminarea acestor dezavantaje am oferit o soluție proprie care se bazează pe o metodă nouă de organizare a subunității logaritmice, pe care am denumit-o **“metoda logaritmului mascat”**, [44]. La elaborarea acestei metode am avut în vedere ca prin aplicarea ei să se atingă și următoarele scopuri:

- noua structură să nu fie organizată pe mai mult de 6 niveluri pipeline,
- să nu crească semnificativ numărul de porți logice utilizate pe fiecare nivel,

- să nu crească semnificativ capacitatea latch-urilor dintre nivelurile structurii pipeline,

- timpul de propagare pe nivelul critic al structurii să fie cât mai aproape de limita fizic realizabilă prin implementarea algoritmului prezentat și anume aceea a timpului de acces la memoriile ROM integrate pe cip,

- timpii de propagare pe nivelurile pipeline să fie cât mai apropiați.

Prin această nouă abordare, termenii y_A , Ey_A , respectiv y_B , Ey_B care se însumează cu produsele $\Delta Ey_A \times y_{2A}$, respectiv $\Delta Ey_B \times y_{2B}$ (conform relației 2.6), vor fi considerați ca și cum ar fi produse parțiale inițiale ale produselor menționate. Acești termeni vor fi astfel incluși în arborele Wallace de reducere a numărului de produse parțiale, pe lângă cele 12 produse parțiale inițiale ale fiecărui produs. În felul acesta au fost practic eliminate sumatoarele finale ale multiplicatoarelor, iar cei doi termeni rezultați în vârful inferior al arborelui, în secțiunile “sumă” și “transport”, vor conține implicit și valoarea produsului, distribuită între cele două secțiuni. Valorile propriu-zise ale celor două produse $\Delta Ey \times y_2$ ne mai regăsim explicit în structură, devin astfel “mascate”.

Problema care mai trebuie rezolvată este aceea a situației în care în relația 2.6 termenul $\Delta Ey \times y_2$ este negativ și ar fi necesară conversia în complement de doi a acestuia, adică a tuturor celor 12 produse parțiale. Acest lucru se întâmplă începând de la adresa 907 și până la adresa 2047 a memoriilor ROMA și ROMA', situație care corespunde sectorului coborât al graficului din figura 2.1.a. Generarea complementului de doi pentru fiecare produs parțial este însă foarte complexă întrucât produsele parțiale au ponderi diferite și conduce la situația din figura 2.15.b. mult mai dificil de implementat decât în cazul când produsul este pozitiv, ca în figura 2.15.a. S-a considerat “p” o reprezentare generică pentru fiecare bit al produsului parțial de 12 biți, indiferent de valoarea sa, iar “q” valoarea bitului “p”, negată. Se constată că pentru a se putea implementa și cazul

b este necesară extinderea la mult mai mulți biți a blocurilor ce reduc numărul de produse parțiale, iar selecția de 1 sau 0 se face pentru un număr foarte mare de intrări în aceste blocuri.

...00000ppppppppppppp
 ...0000ppppppppppppp
 ...000pppppppppppppp
 ...00ppppppppppppppp
 ...0pppppppppppppppp
 ...ppppppppppppppppp

Fig. 2.15.a

Porțiune din aria de înmulțire
 în cazul $\Delta E y \times y_2 > 0$

...11111qqqqqqqqqqqqq
 ...1111qqqqqqqqqqqqq1
 ...111qqqqqqqqqqqqq1
 ...11qqqqqqqqqqqqq1
 ...1qqqqqqqqqqqqq1
 ...1qqqqqqqqqqqqq1
1

Fig.2.15.b

Porțiune din aria de înmulțire
 în cazul $\Delta E y \times y_2 < 0$

În figura 2.15.b bitul de valoare logică 1 de sub fiecare bit q de tip LSB este cel necesar pentru conversia în complement de doi a fiecărui produs parțial.

Am realizat evitarea acestui grav neajuns printr-un artificiu, care permite eliminarea completă a cazurilor în care produsul $\Delta E y \times y_2$ trebuie scăzut. După cum se observă în figura 2.16, se poate scrie următoarea relație:

$$E y_{(n)} - \Delta E y_{(n)} \times y_2 = E y_{(n+1)} + \Delta E y_{(n)} \times (\bar{y}_2 + 1) = E y_{(n+1)} + \Delta E y_{(n)} \times \bar{y}_2 + \Delta E y_{(n)} \quad (2.19)$$

Implementarea acestei relații conduce la aranjarea produselor parțiale așa cum se prezintă în figura 2.17. A fost folosită aceeași reprezentare generică, cu "q_y" pentru biții "p_y" ai lui y₂ complementați, respectiv cu "p_d" pentru biții lui $\Delta E y$.

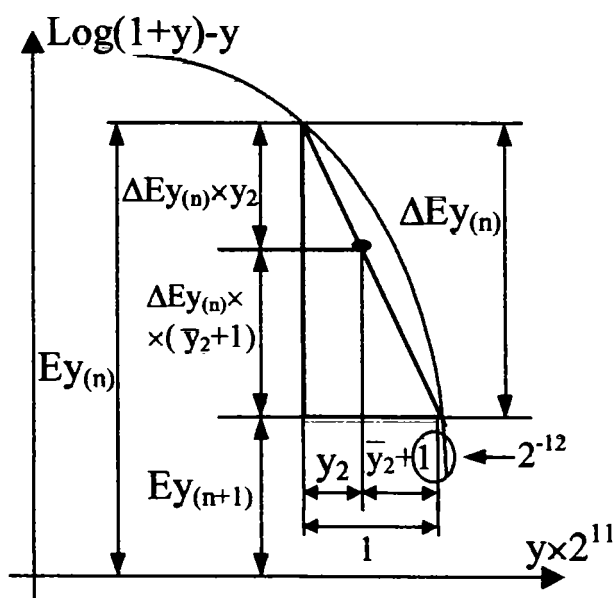


Fig. 2.16

Segment obținut prin interpolare liniară
 între două valori memorate consecutive
 Ey, de pantă negativă

...0 0 0 0 0 p_dp_dp_dp_dp_dp_dp_dp_dp_dp_dp_dp_dp_dp_dp_d
 ...0 0 0 0 0 q_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_y
 ...0 0 0 0 q_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_y
 ...0 0 0 q_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_y
 ...0 0 q_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_y
 ...0 q_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_y
 ...q_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_yq_y

Fig. 2.17

Porțiune din aria de înmulțire în cazul
 implementării relației 2.19

Așadar se poate obține același rezultat al logaritmării dacă se implementează membrul drept al relației 2.19. Începând de la locația de memorie corespunzătoare adresei 907 a lui ROMA, în loc să se memoreze valoarea $E_{y(n)}$, se memorează valoarea $E_{y(n+1)}$, adică cea care ar fi trebuit să se găsească la adresa următoare. În fiecare locație se va memora un bit suplimentar, denumit bit de control, de valoare 0 pentru adresele 0...906 și de valoare 1 pentru adresele 907...2047. Dacă acest bit are valoarea 1 atunci generarea produselor parțiale se va face cu y_2 având biții inversați și se va mai adăuga un produs parțial de pondere egală cu a celui mai puțin semnificativ produs parțial, de valoare $\Delta E_{y(n)}$. Dacă bitul de control are valoarea 0 atunci generarea produselor parțiale se va face cu y_2 neinversat iar biții "p_d" din primul produs parțial din figura 2.17 vor fi toți egali cu 0.

Arborele Wallace va avea astfel drept intrări 15 pseudo-produse parțiale inițiale și va furniza două cuvinte: "sumă" și "transport", așa cum s-a văzut în figura 2.5. Dacă în această fază s-ar efectua adunarea propriu-zisă a acestora s-ar obține valoarea logaritmului mantisei fiecărui operand aplicat la intrarea celor două circuite de logaritmare cu funcționare paralelă. În continuare, cele două numere fracționare obținute ar fi concatenate la exponenții celor doi operanzi, pentru a fi obținuți logaritmi operanzilor. În final cei doi logaritmi s-ar aplica la ALU, pentru a fi adunați sau scăzuți conform relațiilor 2.10 și 2.11. Totuși și în acest caz vor avea loc două adunări propriu-zise consecutive, ceea ce încetinește producerea rezultatului.

Pentru evitarea acestei situații și în acest caz se consideră cele două perechi de cuvinte "sumă" și "transport" ca și pseudo-produse parțiale astfel că ele vor fi introduse într-un nou bloc de reducere. Acesta va furniza în final două cuvinte, "sumă" și "transport" finale, ce vor putea fi adunate apoi cu un sumator rapid. Acest ultim bloc de reducere a numărului de pseudo-produse parțiale va fi inclus în ALU întrucât asupra lui trebuie să acționeze o logică de comandă care să permită implementarea atât a operației de adunare cât și a operației de scădere.

În felul acesta au fost practic eliminate sumatoarele (din blocurile Sumator CSA+CLA, fig. 2.14) care furnizează valorile explicite ale logaritmilor mantiselor celor doi operanzi. Aceste valori devin astfel mascate, ele sunt conținute implicit în cele 4 pseudo-produse parțiale și concură la obținerea rezultatului final. Din acest motiv, această metodă de operare a fost denumită "metoda logaritmului mascat".

Implementarea acestei metode conduce la generarea unui arbore Wallace care are la bază 30 de pseudo-produse parțiale și al cărui vârf se găsește în ALU. Pentru construcția arborelui au fost utilizate drept blocuri de reducere a numărului de pseudo-produse parțiale, blocuri de compresori 4:2 (v. paragraful 2.2.3.2) și blocuri CSA de sumatoare complete 3:2, în așa fel încât să se obțină cea mai eficientă structură, în care să nu existe pseudo-produse parțiale intermediare neoperate pe nici un etaj al arborelui. Structura obținută este prezentată în figura 2.18, iar alinierea diverselor pseudo-produse parțiale inițiale și intermediare de la intrările și ieșirile blocurilor de reducere, în figurile 2.19 și 2.20.

Concepția subunității logaritmice a unui procesor aritmetic hibrid

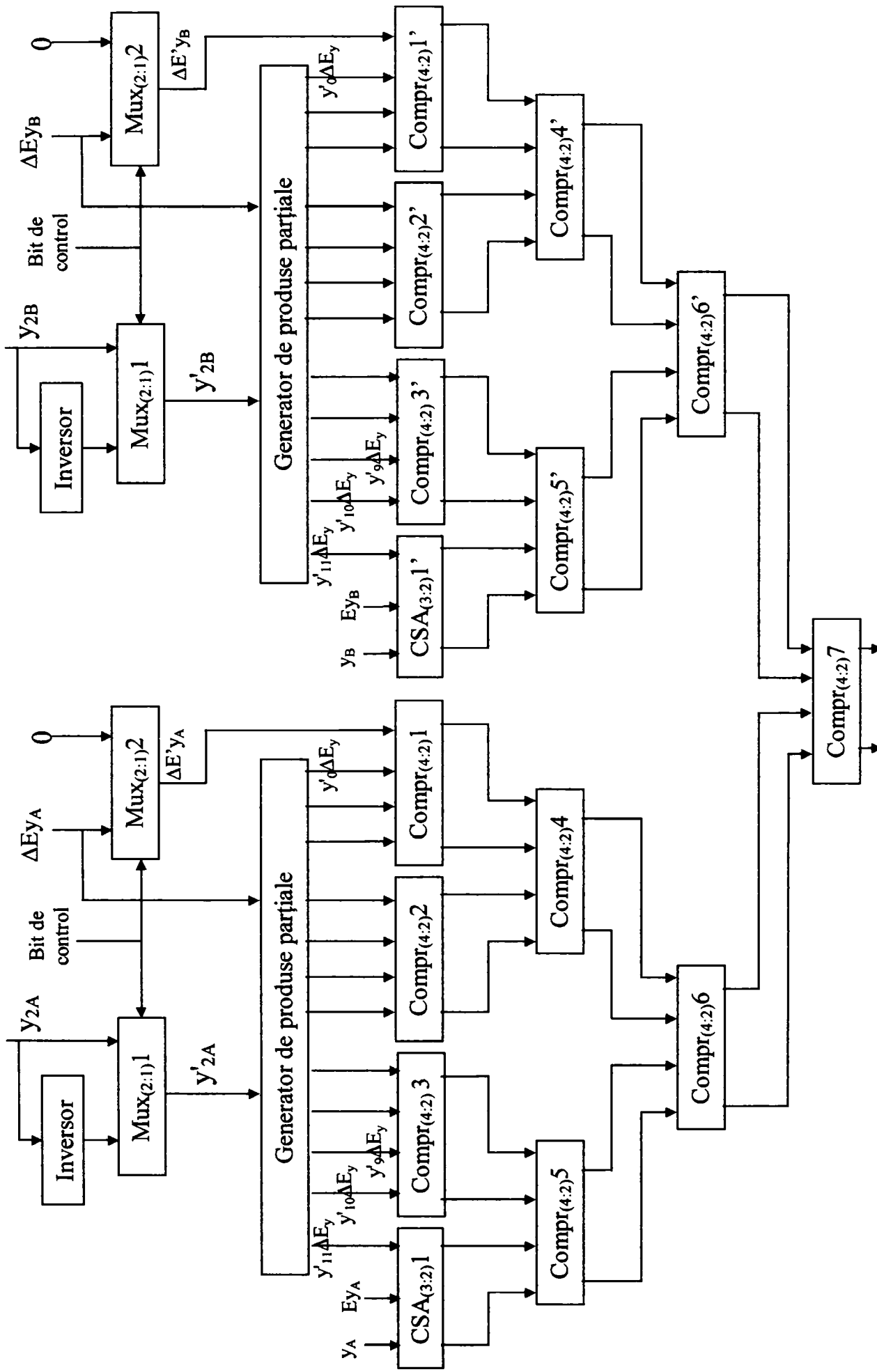


Fig.2.18 Structura arborelui Wallace de reducere a numărului de produse parțiale

În figurile 2.19 și 2.20 a fost folosită notația generică “x” pentru a desemna biții pseudo-produselor parțiale, indiferent de valoarea lor. Blocul CSA grupează termenii y și ΔEy împreună cu cel mai semnificativ produs parțial al produsului $y_2 \times \Delta Ey$, în scopul evitării creșterii exagerate a dimensiunii acestui bloc (și deci și a blocurilor conectate în aval de acesta) și a magistralelor de legătură.

Un bloc de compresori 4:2 va conține în realitate, la anumite ranguri, nu doar compresori 4:2, ci și sumatoare complete 3:2, respectiv sumatoare incomplete 2:2 deoarece, așa cum se observă în figurile 2.19 și 2.20, nu sunt folosite întotdeauna toate cele 4 intrări de același rang. Considerarea acestui fapt va conduce la o economie semnificativă a numărului de porți logice utilizate.

Pentru circuitul de antilogaritmare, procedura aplicată este identică, cuvintele “sumă” și “transport” finale fiind obținute după doar 3 etaje de compresori. Ele sunt apoi însumate cu un sumator rapid pentru obținerea mantisei rezultatului final de la ieșirea subunității logaritmice. Și aici vor fi luate măsuri pentru evitarea operației de scădere a termenului $\Delta Ey \times y_2$, dar se ține cont și de observațiile făcute în paragraful 2.2.2.3. Acestea sunt legate de faptul că în relația 2.9 și termenul Ey trebuie scăzut, iar produsul $\Delta Ey \times y_2$ se scade pe porțiunea crescătoare a curbei din figura 2.1.b și se adună pe porțiunea descrescătoare. Deoarece relația 2.19 nu mai poate fi folosită, suma celor doi termeni negativi din relația 2.9 se scrie astfel:

$$\begin{aligned} -Ey_{(n)} - \Delta Ey_{(n)} \times y_2 &= -(Ey_{(n)} + \Delta Ey_{(n)} \times y_2) = -[Ey_{(n)} + \Delta Ey_{(n)} - \Delta Ey_{(n)} \times (\bar{y}_2 + 1)] = \\ &= -(Ey_{(n)} + \Delta Ey_{(n)}) + \Delta Ey_{(n)} \times \bar{y}_2 + \Delta Ey_{(n)} = -Ey_{(n+1)} + \Delta Ey_{(n)} \times \bar{y}_2 + \Delta Ey_{(n)} \end{aligned} \quad (2.20)$$

Așadar în locațiile memoriei ROMC de la adresa 0 la adresa 1082 se va memora complementul de doi al cantităților $Ey_{(n+1)}$, ce ar fi trebuit să se găsească la adresa următoare, precum și valoarea 1 pentru bitul de control; de la adresa 1083 la adresa 2047 (când $\Delta Ey \times y_2$ este pozitiv) se va memora direct complementul de doi al valorilor $Ey_{(n)}$, precum și valoarea 0 pentru bitul de control.

2.2.3.2 Prezentarea noului tip de compresor 4:2 folosit în proiectarea subunității logaritmice

Compresorul 4:2 a fost introdus în literatura de specialitate de către Weinberger în 1981. Prima implementare pe cip ce făcea apel la acest circuit apare însă abia în 1991, prin multiplicatorul 54×54 de biți conceput de J. Mori [67], posibil de realizat în urma creșterii densității de integrare. De atunci au apărut și alte variante îmbunătățite de compresor 4:2, care vor fi comparate la finele paragrafului cu varianta proprie a acestui circuit, denumită compresor 4:2 L.

Denumirea generalizată de “compresor 4:2” nu reflectă însă întocmai caracteristicile circuitului deoarece suma celor 4 biți de intrare nu poate fi reprezentată în toate cazurile posibile, cu ajutorul celor doi biți de ieșire. Circuitul furnizează în plus un transport pe o linie suplimentară de ieșire și are, în

consecință, nevoie de o linie pentru introducerea transportului de intrare. Există așadar o cale de propagare pe verticală, “4:2” și o cale de propagare pe orizontală “1:1”. Totuși propagarea pe orizontală a transportului este practic limitată la un singur bit, întrucât transportul de ieșire este generat astfel încât să nu depindă de transportul de intrare. Astfel, la scara întregului multiplicator, utilizarea acestui circuit permite deplasarea căii critice de propagare din centrul ariei de înmulțire (unde este gradul maxim de suprapunere a produselor parțiale) spre capătul din stânga (unde acest grad de suprapunere devine din ce în ce mai mic).

În plus, combinația dintre blocuri CSA de sumatoare complete 3:2 și blocuri de compresori 4:2, permite construcția eficientă a arborelui Wallace, prin eliminarea cazurilor în care una sau două produse parțiale inițiale sau intermediare rămase negrupate traversează un etaj al arborelui fără a fi operate [97]. Numărul de etaje al arborelui Wallace devine astfel mai mic, respectiv se scurtează calea critică de propagare prin acesta.

Principial, un compresor 4:2 se obține prin interconectarea a două sumatoare complete 3:2, a căror structură a fost prezentată în figura 1.1 și care prezintă un timp de propagare echivalent a 4 porți SAU-exclusiv cascade.

O îmbunătățire majoră a structurii compresorului 4:2 a fost adusă de D. Villegier în [106] și [72] prin exploatarea intrărilor și ieșirilor “rapide” ale sumatoarelor complete 3:2, în așa fel încât timpul maxim de propagare devine echivalentul a doar 3 porți SAU-exclusiv cascade. Această structură optimizată este prezentată în figura 2.21.

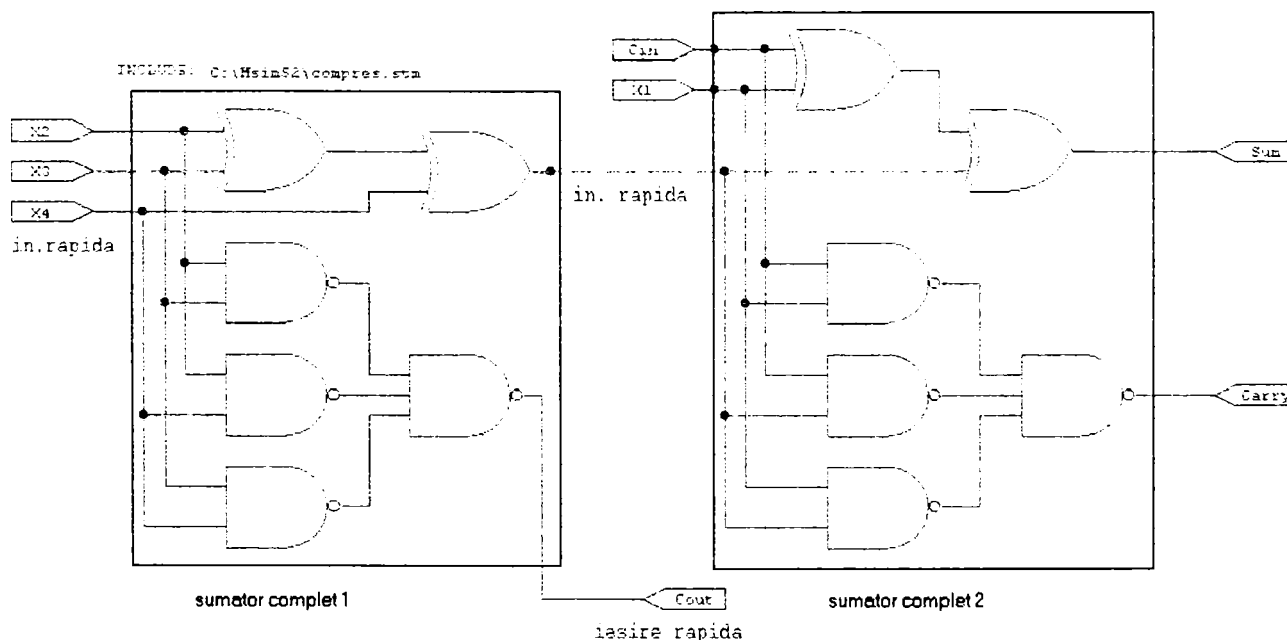


Fig.2.21 Structura unui compresor 4:2 optimizat [72]

Pentru acest compresor, cei 4 biți de intrare sunt x1, x2, x3 și x4, iar cei 2 biți de ieșire sunt S_{um} și C_{arry}. Transportul de intrare este C_{in} iar transportul de ieșire este C_{out}. Biții C_{arry} și C_{out} au pondere dublă față de bitul S_{um}. Bitul Sum este obținut prin implementarea funcției SAU-exclusiv între cele 5 intrări ale circuitului. În presupunerea că timpul de propagare a transportului printr-o poartă

SAU-exclusiv este de două ori mai mare decât printr-o poartă ȘI-NU, se poate observa din figura 2.21 că, indiferent de calea de propagare a semnalelor de intrare, inclusiv a lui C_{in} generat de compresorul din amonte, timpul maxim de propagare nu depășește echivalentul a șase porți logice simple sau trei porți SAU-exclusiv.

Se mai observă că ieșirea C_{out} nu depinde de intrarea C_{in} , deci propagarea transportului pe orizontală în blocul de compresori din figura 2.22 este limitată la un singur bit (o poziție) spre stânga. Timpul de propagare pe verticală al întregului bloc este de fapt timpul de propagare pe verticală printr-un singur compresor. Aceasta înseamnă evident că înjumătățirea numărului de produse parțiale se face într-un timp egal cu timpul de propagare pe calea critică prin compresorul 4:2. Din acest motiv, propagarea rapidă a semnalelor prin acest circuit este esențială și determină practic viteza de propagare prin întreg arborele Wallace.

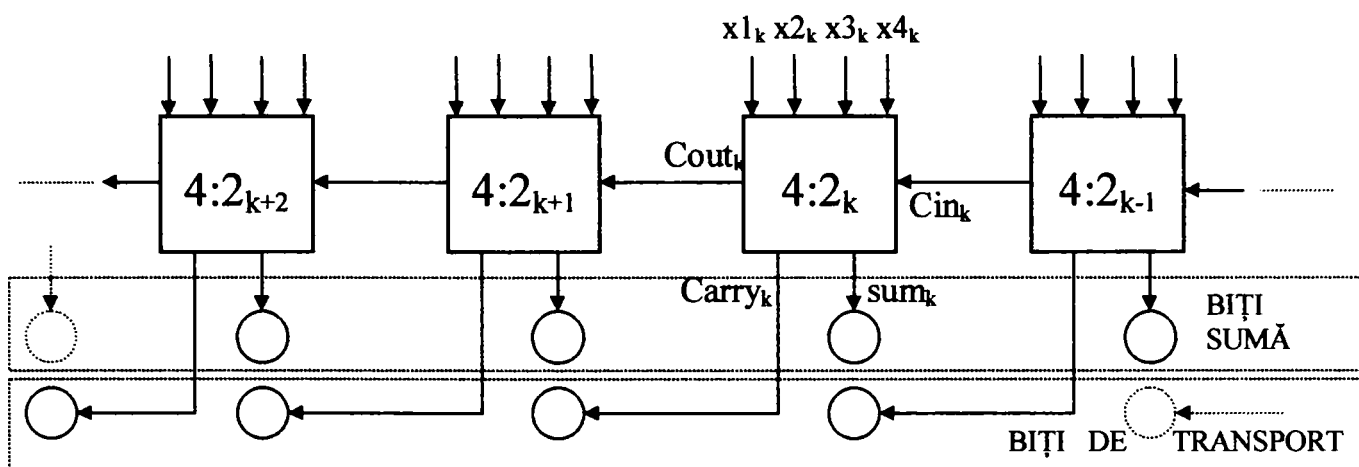


Fig. 2.22 Secțiune printr-un bloc de compresori 4:2

Compresorul 4:2 optimizat descris în lucrarea [72] și prezentat în figura 2.21 manifestă însă o serie de limitări care afectează viteza de propagare a transportului prin el:

- porțile ȘI-NU cu trei intrări, în tehnologii submicronice avansate ce necesită tensiuni de alimentare mai mici, nu mai pot fi implementate sau se implementează prin porți ȘI-NU cu două intrări, prin adăugarea încă a unui nivel logic. În acest fel calea critică de propagare prin circuit devine cea care include ieșirea rapidă C_{out} a compresorului actual și intrarea C_{in} și ieșirea C_{arry} a compresorului din aval. Această cale devine deci mai lentă decât calea de propagare pentru "x1" și "x2" (fig. 2.21) prin trei porți SAU-exclusiv succesive; chiar și în tehnologii care permit integrarea porții ȘI-NU cu trei intrări, timpul de propagare prin aceasta este mai mare decât printr-o poartă cu două intrări.

- utilizarea unor porți SAU-exclusiv neconvenționale [67], mai rapide decât cele clasice, având $t_{P(SAU-excl)} < 2 \times t_{P(ȘI-NU)}$ nu va aduce avantajele scontate din cauză că propagarea prin trei porți SAU-exclusiv cascade nu mai reprezintă o cale critică în compresor.

- există multe cazuri în care încărcarea capacitivă a unei ieșiri logice este considerată excesivă, atunci când ea apare pe o cale critică de propagare. Există

astfel în structura compresorului cazuri în care o ieșire comandă câte trei intrări, adică minim șase grile MOS. La capacitatea intrinsecă a acestora se adaugă și capacitatea parazită a celor trei linii de semnal.

Eliminarea acestor dezavantaje s-a făcut prin conceperea unei variante proprii de compresor 4:2, [38], care așa cum s-a dovedit prin simularea pe layout, este mai rapidă cu 10 % decât varianta din [72]. Ținând cont că în subunitatea logaritmă există 7 etaje de compresori 4:2 (4 etaje în circuitul de logaritmare și ALU și 3 în circuitul de antilogaritmare), aceasta duce în ansamblu la o economie de 0,7 ns (corespunzător implementării în tehnologie CMOS 0,5 microni).

În conceperea variantei proprii am plecat de la ideea găsirii unei definiții generale pentru compresorul 4:2, ale cărui intrări și ieșiri sunt deja definite. Astfel acest circuit trebuie să asigure următoarele funcții:

a) ieșirea S_{UM} să fie generată de relația 2.21:

$$S_{um} = x1 \oplus x2 \oplus x3 \oplus x4 \oplus C_{in} \quad (2.21)$$

b) Cele două linii pe care apar transporturile de ieșire C_{arry} și C_{out} să fie de rang imediat superior lui S_{um} (să aibă pondere dublă), astfel încât valoarea maximă a sumei celor 5 biți de intrare, egală cu 5, să poată fi reprezentată. Astfel va fi îndeplinită relația 2.22:

$$x1+x2+x3+x4+C_{in} = C_{arry} + C_{out} + S_{um} \quad (2.22)$$

Într-adevăr, atunci când toți biții de intrare sunt 1 se obține $5=10_2+10_2+1$.

c) Una din liniile de transport de ieșire să poată fi generată de funcții logice simple, care să includă un număr minim de niveluri logice și să nu depindă de transportul de intrare C_{in} . Ea se numește ieșire rapidă și va asigura transportul de intrare pentru compresorul de rang imediat superior. În felul acesta propagarea pe orizontală a transportului se va limita la un singur bit și nu va influența viteza de generare a rezultatului la ieșirea unui bloc de compresori 4:2 cascadați.

d) Toate combinațiile posibile de biți de intrare pentru care suma lor este 0 sau 1, trebuie să conducă la $C_{arry} = C_{out} = 0$.

e) Toate combinațiile posibile de biți de intrare pentru care suma lor este 2 sau 3, trebuie să conducă fie la $C_{arry} = 1$ și $C_{out} = 0$, fie la $C_{arry} = 0$ și $C_{out} = 1$. Cu alte cuvinte, toate aceste combinații posibile care produc fie $C_{arry} = 1$, fie $C_{out} = 1$ trebuie să se găsească în două mulțimi disjuncte.

f) Toate combinațiile posibile de biți de intrare pentru care suma lor este 4 sau 5, trebuie să conducă la $C_{arry} = C_{out} = 1$.

Având în vedere toate aceste funcții, văzute ca și cerințe obligatorii, dar și dezavantajele compresorului 4:2 de referință, a rezultat structura eficientă a compresorului 4:2 din figura 2.23, denumit compresor 4:2 L.

Datorită numărului mare de intrări, pentru generarea bitului S_{um} conform relației 2.20, am adoptat o structură arborescentă de porți SAU-exclusiv, care grupează intrările $x1$, $x2$, $x3$, $x4$ și ulterior C_{in} . În continuare, pe baza grupării primilor patru biți de intrare, a fost implementat C_{out} prin funcția logică simplă:

$$C_{out} = x_1 x_2 + x_3 x_4 \quad (2.23)$$

Se obține astfel generarea ieșirii rapide “ C_{out} ” prin parcurgerea a doar două nivele logice, realizată exclusiv cu porți intrinseci cu doar două intrări, aplicând relațiile De Morgan. Nici una din liniile de intrare nu mai este astfel încărcată decât cu maxim două intrări logice. Se observă că relația 2.23 îndeplinește cerințele d) și f). Pentru obținerea ieșirii “ C_{arry} ” a fost implementată relația 2.24:

$$C_{arry} = (x_1 \oplus x_2)(x_3 \oplus x_4) + C_{in}(x_1 \oplus x_2) + C_{in}(x_3 \oplus x_4) + x_1 x_2 x_3 x_4 \quad (2.24)$$

Primii trei termeni ai sumei logice asigură îndeplinirea cerinței e) relativ la C_{out} , iar termenul al patrulea asigură îndeplinirea cerinței f) din lista enunțată. De asemenea se observă că cerința d) este îndeplinită de către toți cei patru termeni.

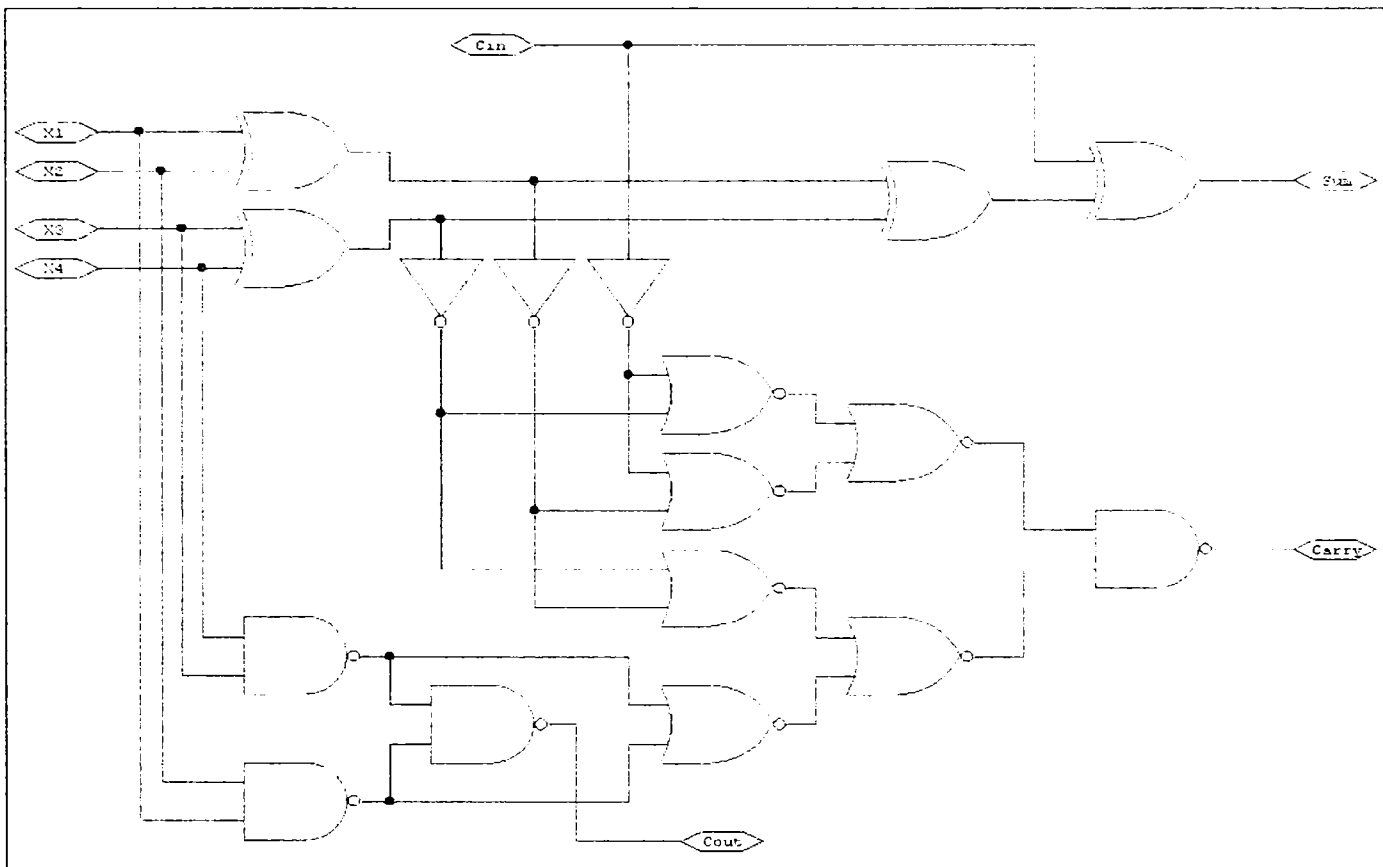


Fig. 2.23 Structura compresorului 4:2 L

Utilizarea celor trei inversoare din figură a permis eliminarea completă a dezavantajului încărcării suplimentare capacitive a anumitor ieșiri logice (o ieșire nu comandă mai mult de două intrări) și în același timp a permis aplicarea relațiilor De Morgan, pentru a utiliza exclusiv porți intrinseci în tehnologie CMOS. O comparație de viteză între varianta proprie de compresor 4:2 și câteva variante recente, apărute în literatura de specialitate, pentru o tehnologie CMOS de 0,5 microni, este ilustrată în tabelul 2.3:

Tab. 2.3

Tip compresor	[67] ₁₉₉₁	[70] ₁₉₉₅	[72] ₁₉₉₆	[101] ₁₉₉₉	Compres. 4:2 L
Timp de propagare	1,4 ns	1,4 ns	1,3 ns	1,2 ns	1,1 ns

Așa cum se observă, compresorul 4:2 L este superior în viteză celorlalte variante.

2.2.3.3 Concepția blocului ALU în cazul aplicării metodei logaritmului mascat

Datorită faptului că în blocul ALU trebuie implementată pe lângă operația de adunare și operația de scădere, ce implică realizarea conversiei în complement de doi a scăzătorului, în mod obligatoriu blocul de compresori 4:2 L, Compr_(4:2)⁷ din figura 2.18 trebuie inclus în ALU. El va efectua reducția celor 4 pseudo-produse parțiale de intrare după ce logica de comandă a circuitului sumator/scăzător a intervenit asupra scăzătorului.

Scăzătorul este reprezentat de două pseudo-produse parțiale, a căror însumare nu se mai execută, tocmai în scopul creșterii vitezei de propagare a transportului prin subunitatea logaritmice. Aceasta înseamnă că ambii termeni trebuie convertiți în cod complement de doi. Evitarea reconversiei din cod complement de doi în cod semn-mărime a rezultatului ce va fi obținut, în cazul când acesta este negativ, se face după aceeași metodă ca și în paragraful 2.2.2.4, ce uzează de topologia de concepție proprie a circuitului sumator/scăzător.

Se notează cu A_1, A_2 , respectiv B_1, B_2 cele 4 pseudo-produse parțiale, caz în care $A=A_1+A_2$ reprezintă descăzutul, iar $B=B_1+B_2$ reprezintă scăzătorul, în cazul în care se efectuează o operație de scădere. Considerând convenția de reprezentare a lui B negat (data obținută prin inversarea tuturor biților lui B) ca fiind “/B”, se pot scrie cei doi termeni care se calculează simultan în circuitul sumator/scăzător:

$$A - B = A - B_1 - B_2 = A + /B_1 + /B_2 + 2 \quad (2.25)$$

$$/(B - A) = A + /B_1 + /B_2 + 1 \quad (2.26)$$

Relația 2.26 se verifică în felul următor:

- în ecuația 2.27 de mai jos,

$$A - B = - (B - A) = / (B - A) + 1 \quad (2.27)$$

se înlocuiește termenul “/(B - A)” cu valoarea dată de relația 2.26.

- se obține într-adevăr $A - B = (A + /B_1 + /B_2 + 1) + 1 = A + /B_1 + /B_2 + 2$, adică exact valoarea corectă din relația 2.25.

Principiul circuitului sumator/scăzător este următorul:

- dacă A-B rezultă pozitiv, el este selectat la ieșire, iar bitul de semn al rezultatului va fi 0;

- dacă A-B rezultă negativ înseamnă că B-A este pozitiv și la ieșire se selectează termenul /(B-A) negat, adică B-A. Bitul de semn al rezultatului va fi 1.

Pentru a menține similitudinea cu structura din figura 2.9, relațiile 2.25 și 2.26 se rescriu astfel:

$$A - B = (A_1 + A_2 + /B_1 + /B_2 + 1) + 1 \quad (2.28)$$

$$/(B - A) = (A_1 + A_2 + /B_1 + /B_2 + 1) + 0 \quad (2.29)$$

Cu alte cuvinte se menține situația de transport inițial 1, respectiv transport inițial 0 la cele două sumatoare care lucrează în paralel, condiționat de asigurarea unui transport de intrare egal cu 1 la blocul de compresori 4:2 L, în cazul efectuării

operației de scădere. Evident acest transport inițial va fi 0 în cazul efectuării operației de adunare. Transportul de intrare se aplică pe intrarea C_{in} nefolosită a compresorului 4:2 L de rangul cel mai mic.

În continuare, dimensiunea blocului $Compr_{(4,2)7}$, inclus în ALU, se va suplimenta cu 9 biți la stânga, pentru concatenarea exponenților pozitivi (cu deplasamentul de 127 inclus, conform algoritmului modificat) ce reprezintă partea întreagă a logaritmului cu deplasament a celor doi operanzi de intrare. Părțile fracționare ale logaritmilor celor doi operanzi sunt "mascate", fiind distribuite între elementele A_1 și A_2 , respectiv B_1 și B_2 ale celor două perechi de pseudo-produse parțiale. Concatenarea exponenților se va face la termenii A_1 și B_1 , obținându-se pseudo-produsele parțiale finale \underline{A}_1 și \underline{B}_1 , în timp ce cele 9 poziții de biți de parte întreagă corespunzătoare la A_2 și B_2 de pondere fracționară, se vor completa cu zerouri. Se vor obține astfel celelalte două pseudo-produsele parțiale finale \underline{A}_2 și \underline{B}_2 .

Schema bloc a noului sumator/scăzător, conceput în scopul implementării complete a metodei logaritmului mascat, este prezentată în figura 2.24. Linia S_{OP} are valoarea logică 0 în cazul adunării și 1 în cazul scăderii.

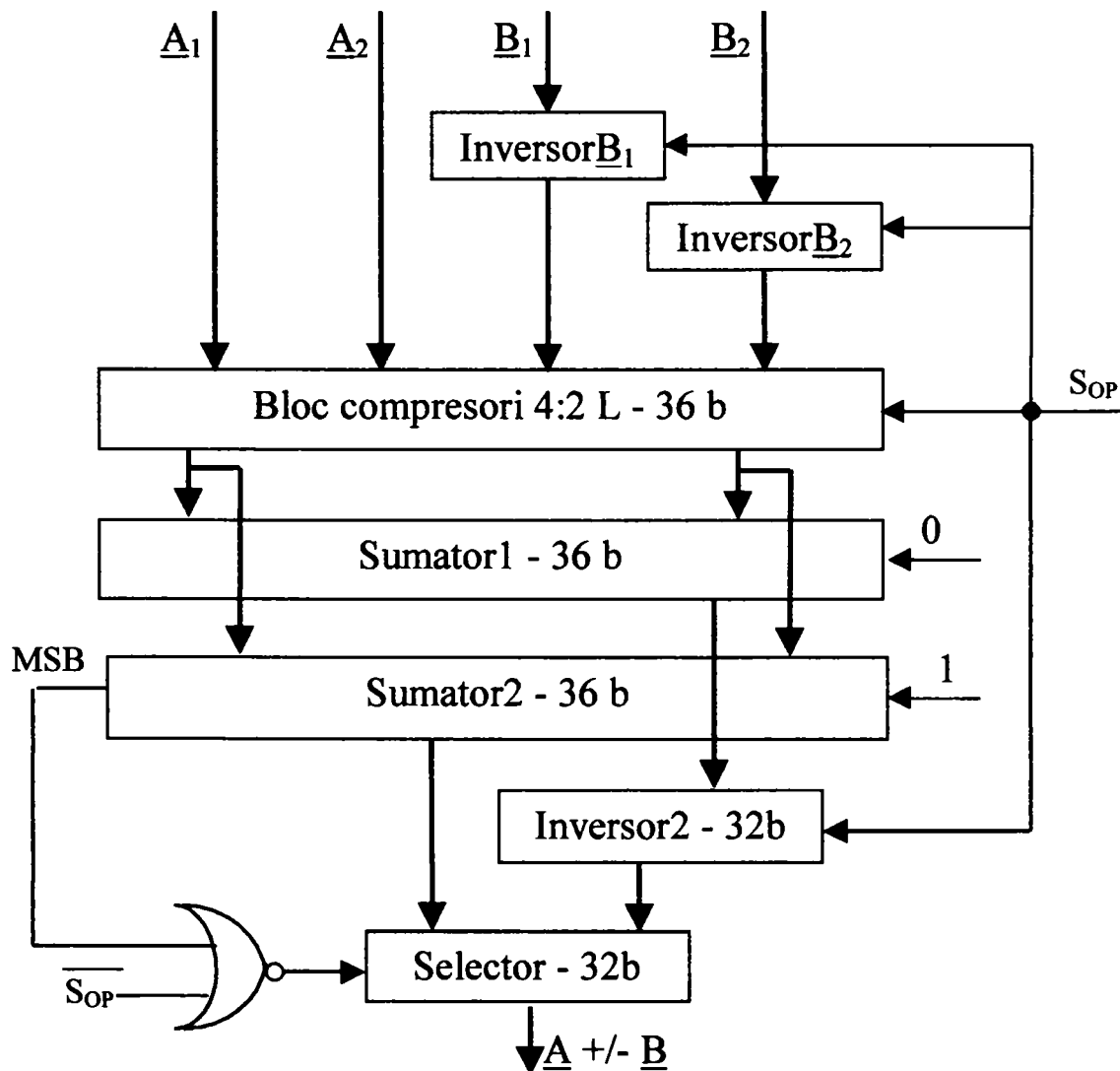


Fig. 2.24 Schema bloc a noului circuit sumator/scăzător

Blocurile Inversor \underline{B}_1 , Inversor \underline{B}_2 și Inversor2 vor fi transparente atunci când se efectuează o operație de adunare și vor inversa biții datelor de la intrări atunci când se efectuează o operație de scădere. Blocul Selector va selecta rezultatul de la Sumator1 neinvertat, în cazul adunării, iar în cazul scăderii va selecta rezultatul de la Sumator2 sau de la Sumator1 invertat, în funcție de valoarea bitului MSB a blocului Sumator2. Detaliile de funcționare sunt similare cu cele prezentate în paragraful 2.2.2.4.

Cele două sumatoare funcționează pe 36 de biți, fiind implementate prin același tip de circuit performant de concepție proprie, care a fost prezentat în figura 2.12. Diferența constă în aceea că în structura de sumator cu selecția transportului mai apare o secțiune suplimentară CLA pe 4 biți, în poziția cea mai semnificativă. Cei mai puțin semnificativi 4 biți de la ieșirile celor două sumatoare vor fi pierduți, astfel că la ieșirea circuitului se regăsește dimensiunea formatului simplă precizie plus un bit în poziția MSB, care evită apariția supradepășirii datorată cumulării deplasamentelor în cazul înmulțirii.

Justificarea dimensiunii de 36 de biți a sumatoarelor și a blocului de compresori va fi făcută în paragraful 2.2.4.

Schemele detaliate ale întregului circuit, cât și ale arborelui Wallace din amonte sunt prezentate în Anexa 3.

Timpul de propagare maxim prin arborele Wallace ce conține trei etaje de compresori 4:2 L este $3 \times 1,1$ ns, la care se mai adaugă 0,7 ns pentru generarea produselor parțiale ceea ce conduce la un total de 4 ns, pentru tehnologie CMOS 0,5 micrometri.

Timpul de propagare maxim pentru circuitul sumator/scăzător în cazul în care sumatoarele sunt de 32 de biți, este obținut prin cumularea valorii din tabelul 2.2 din paragraful 2.2.2.4, adică 2,8 ns, cu timpul de propagare prin blocul de compresori 4:2 L, de 1,1 ns, rezultând 3,9 ns. În cazul de față se mai adaugă maxim 0,2 ns necesare pentru multiplicarea cu 2 a capacității de comandă a ultimei linii de selecție din sumatorul de 36 de biți. Evaluarea cazului celui mai defavorabil, pentru a se efectua simularea, este dificilă datorită combinației dintre cele două tipuri de circuite, care pentru o valoare precizată a intrărilor pot să nu se găsească simultan în cazul cel mai defavorabil. De aceea, prin simulare se pot obține timpi mai mici de 4 ns.

Pentru sumatorul final pe 27 de biți de la capătul arborelui Wallace din circuitul de antilogaritmă, de același tip cu sumatorul din figura 2.11, a fost găsit un timp de propagare în cazul cel mai defavorabil de 3,6 ns.

2.2.3.4 Concluzii la paragraful 2.2.3

Utilizarea metodei logaritmului mascat în implementarea subunității logaritmice de calcul a condus la obținerea următorului timp total de producere a rezultatului (au fost luați în considerare timpii de propagare maximi ai tuturor blocurilor aflate pe calea critică de propagare):

$$\begin{aligned} t_{\text{tot}} &= t_{\text{ROMA}} + t_{\text{arbore.log}} + t_{\text{ALU}} + t_{\text{ROMC}} + t_{\text{arbore.alog}} + t_{\text{SumFinal}} = \\ &= 3,8 \text{ ns} + 4 \text{ ns} + 4 \text{ ns} + 3,8 \text{ ns} + 4 \text{ ns} + 3,6 \text{ ns} = 23,2 \text{ ns} \end{aligned} \quad (2.30)$$

Se observă că în acest caz, spre deosebire de rezultatele obținute în paragraful 2.2.2 care erau comparabile în viteză cu cea mai performantă implementare descrisă în literatura de specialitate de către F. Lai [50], timpii de propagare prin blocurile de pe calea critică de propagare sunt considerabil mai mici și sunt mult mai echilibrați. Acest lucru permite o organizare foarte eficientă într-o structură pipeline cu 6 niveluri, așa cum se observă în figura 2.25. Așa cum se va dovedi prin simulare, există rezerve de $0,2 \div 0,4 \text{ ns}$ relativ la cele 4 ns din fiecare nivel, care compensează timpul de încărcare a latch-urilor pipeline.

Frecvența de clock a sistemului ajunge astfel la aproape 250 MHz, ceea ce reprezintă o creștere de viteză cu 60% față de referința menționată. Toate avantajele legate de viteza de calcul, pe care le aduce un astfel de procesor și evidențiate în paragraful 2.2.2, se multiplică cu un factor de 1,6.

În plus, se pot remarca și alte caracteristici notabile:

- *numărul de porți logice utilizate în întreaga structură a scăzut semnificativ, datorită eliminării a șase sumatoare rapide de minim 18 biți, cu prețul suplimentării cu 4 biți a trei sumatoare rapide;*

- *timpul de propagare prin nivelul critic al structurii pipeline a coborât practic la limita fizică minimă, dată de timpul de acces la memorii;*

- *capacitatea latch-urilor dintre nivelurile pipeline a scăzut cu aproximativ 30 %, fapt datorat structurii arborescente a subunității logaritmice;*

- *o economie suplimentară de suprafață, de aproximativ 25%, s-a realizat la implementarea celor trei înmulțiri de tip $\Delta E_y \times y_2$, posibilitate demonstrată în urma unui calcul al erorilor în cazul efectuării produsului menționat. Acest lucru va fi prezentat în paragraful 2.2.4.*

2.2.4 Studiul erorilor produse la generarea logaritmului și antilogaritmului binar

În urma analizei cu programul Matlab, a erorilor introduse prin utilizarea algoritmului prezentat în paragraful 2.2.1, s-a confirmat valoarea de 3×10^{-7} menționată în [50], ca valoare maximă a erorii de conversie. Aceasta eroare este totuși de circa 2,5 ori mai mare decât eroarea în formatul simplă precizie, care este de $1,19 \times 10^{-7}$. Pentru minimizarea în continuare a erorii de conversie am propus o soluție proprie, care va adăuga valori de corecție pe anumite intervale de adresă ale memoriilor de tip ROMA și ROMC, în locațiile corespunzătoare de memorie. Ținând cont că eroarea în simplă precizie, adică valoarea celui mai puțin semnificativ bit pe care o furnizează o ieșire oarecare a memoriei ROMA sau ROMC, este de $1,2 \times 10^{-7}$, înseamnă că la valorile calculate ale lui E_y pot fi adăugate corecții de ordinul a un LSB, după care sunt memorate direct (ROMA), respectiv se transformă în complement de doi și apoi se memorează (ROMC).

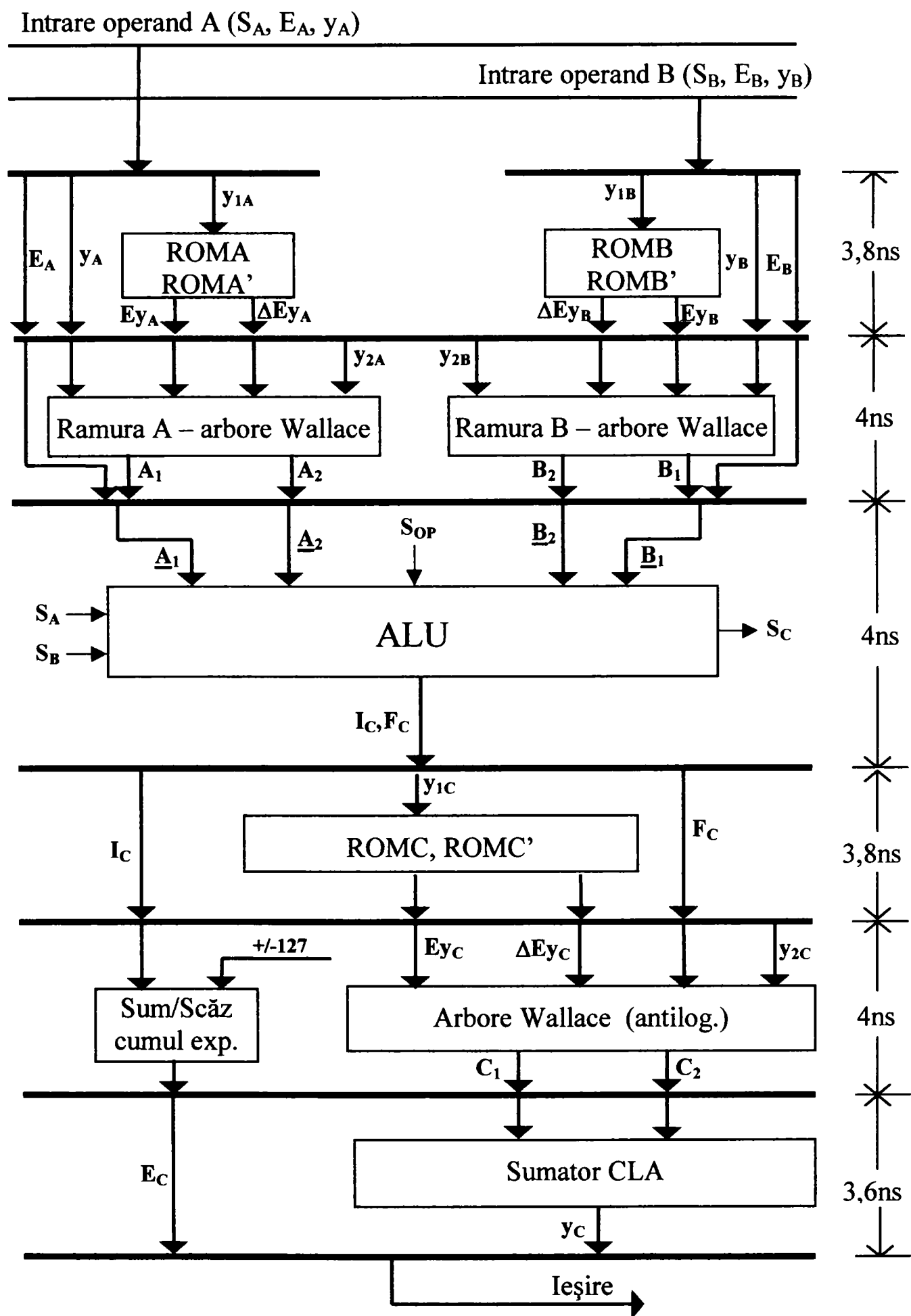


Fig. 2.25 Organizarea noii subunități logaritmice pe niveluri pipeline

Așa cum se vede în figura 2.26, teoretic eroarea totală de conversie se poate reduce la jumătate, dacă la valoarea memorată E_y se adaugă o cantitate egală cu jumătate din eroarea maximă de conversie, adică $1,5 \times 10^{-7}$. Această cantitate nu poate fi reprezentată pe un număr întreg de biți, care ar fi putut să fie 1 sau 2. Soluția de compromis este deci adăugarea unui 1 (LSB) în acele locații pentru care eroarea de conversie este mai mare decât 1 LSB, adică mai mare decât $1,2 \times 10^{-7}$.

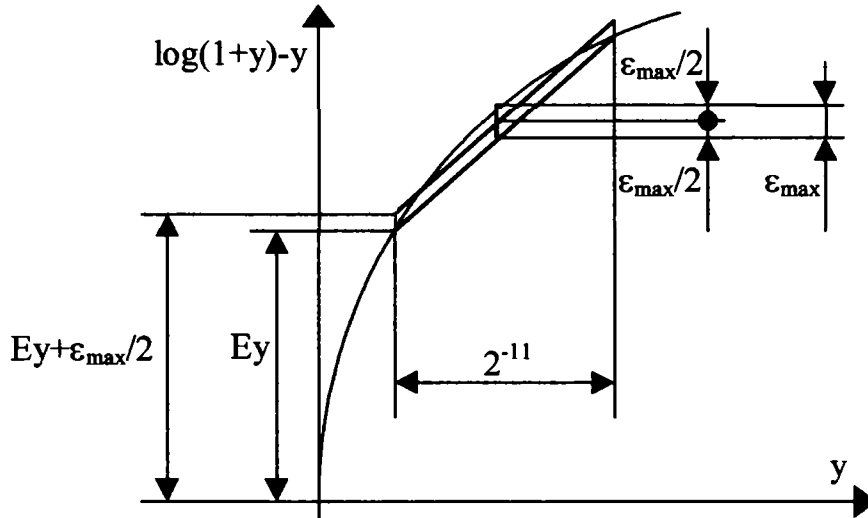


Fig. 2.26 Model teoretic pentru reducerea la jumătate a erorii maxime de conversie

În Anexa 4 sunt prezentate programele Matlab care identifică intervalele de adresă în care trebuie făcută corecția cu 1 bit și furnizează valorile optime E_y , necesare pentru calculul logaritmului și antilogaritmului. Aceste programe vor furniza și valorile ΔE_y necesare pentru cele două conversii de format. În felul acesta, eroarea de conversie, atât la logaritmare cât și la antilogaritmare a scăzut la $1,8 \times 10^{-7}$.

În ce privește calcularea produsului $\Delta E_y \times y_2$, se poate constata din figura 2.19 unde este prezentată aria de înmulțire, că dacă se face calculul erorii de trunchiere, în cazul cel mai defavorabil, când toți biții de rang mai mic sau egal cu “-28” (marcați cu caractere neîngroșate) sunt egali cu 1, se obține valoarea:

$$2 \times 10^{-35} + 3 \times 10^{-34} + 4 \times 10^{-33} + 5 \times 10^{-32} + 6 \times 10^{-31} + 7 \times 10^{-30} + 8 \times 10^{-29} + 9 \times 10^{-28} = 5,96 \times 10^{-8}$$

Această valoare reprezintă jumătate din eroarea de reprezentare în simplă precizie.

Înlăturarea tuturor biților din zona în care ei sunt marcați cu caractere neîngroșate, din figura 2.19 și figura 2.20, adică eliminarea structurilor hard din tot arborele Wallace și ALU care corespund acestei zone, nu introduce eroare suplimentară dacă valorile E_y care se memorează în ROMA se rotunjesc la 23 de biți prin adaos, iar valorile E_y al căror complement de doi se memorează în ROMC, se rotunjesc prin lipsă.

Cap.3 Concepția unei subunități de calcul în virgulă flotantă pentru realizarea operațiilor de adunare și scădere

În acest capitol va fi prezentată o subunitate de calcul în virgulă flotantă, care implementează operațiile de adunare și scădere. Structura proiectată va fi optimizată din punct de vedere al suprafeței ocupate pe cip și al realizării sincronizării cu subunitatea logaritmică. Astfel, ansamblul format din cele două subunități va funcționa în parametri de mare performanță, mai ales în cazul implementării prin microprogramare a unor algoritmi DSP, sau transformări de coordonate pentru aplicații grafice.

Ținând cont că procesorul de calcul descris în această lucrare a fost conceput în primul rând pentru acest gen de aplicații și apoi ca și coprocesor propriu-zis, nu atât de importantă se consideră a fi durata totală a unei operații singulare (ce depinde de numărul maxim de niveluri al structurii pipeline), cât timpul de propagare a transportului pe un nivel al structurii pipeline. Un algoritm DSP poate presupune zeci sau chiar sute de perioade de clock și atunci contează mai puțin faptul că la numărul total de cicli se mai adaugă încă câțiva pentru introducerea și procesarea progresivă a primilor operanzi respectiv pentru recuperarea ultimelor date procesate.

În acest context și în conformitate cu aspectele analizate în paragraful 1.2.3 legate de organizarea unei subunități de adunare/scădere în virgulă flotantă, am adoptat pentru aceasta o structură pipeline pe trei niveluri, care să poată funcționa la o frecvență de clock de 250 MHz, ca și subunitatea logaritmică, atunci când tehnologia utilizată este CMOS 0,5 μm . Deși literatura de specialitate oferă soluții care permit reducerea numărului de niveluri pipeline la două în condițiile menținerii aceleiași frecvențe de lucru a sistemului, prețul plătit prin creșterea exagerată a suprafeței ocupate pe cip de această subunitate de calcul (așa cum s-a menționat în paragraful 1.2.3) nu justifică creșterea puțin semnificativă a vitezei de procesare în cazul unor aplicații DSP.

Datele inițiale, respectiv parametrii între care trebuie să se încadreze subunitatea de adunare/scădere în virgulă flotantă se referă așadar la organizarea sa pe trei niveluri pipeline, iar timpul maxim de propagare al transportului pe fiecare nivel nu trebuie să depășească 4 ns. Astfel, în nivelul 1 se efectuează alinierea mantiselor, în nivelul 2 se operează mantisele aliniate, iar în nivelul 3 se realizează normalizarea rezultatului. Ca și în cazul subunității logaritmice, rotunjirea rezultatului se va face implicit, prin lipsă.

Au fost proiectate și simulate cu ajutorul programului MSim circuite rapide de aliniere a mantisei și de normalizare precum și un sumator/scăzător în virgulă fixă pe 24 de biți. Așa cum se va vedea ulterior, s-au obținut timpi de propagare pe fiecare nivel pipeline sub 3 ns, comparabili sau chiar mai mici decât în cazul unor implementări prezentate în literatura de specialitate. În aceste condiții, pentru a menține sincronizarea cu unitatea logaritmică la un preț plătit cât mai mic, a avut

loc o reorientare în concepția subunității în virgulă flotantă, în sensul optimizării ei din punct de vedere al suprafeței ocupate pe cip. Complexitatea circuitelor necesare pentru implementarea operațiilor de adunare și scădere a fost redusă până la acel punct ce corespunde situației în care durata de procesare a informației binare pe fiecare nivel pipeline încă nu a depășit 4 ns.

3.1 Proiectarea circuitului de aliniere a mantiselor

Se știe că deplasarea la dreapta cu un bit a mantisei implică incrementarea exponentului asociat cu o unitate; respectiv dacă diferența exponenților celor doi operanzi este 1, mantisa asociată celui mai mic exponent trebuie deplasată la dreapta cu un bit. Aceasta înseamnă că în general dacă bitul de ordin i (s_i) din diferența exponenților este 1, mantisa corespunzătoare celui mai mic exponent trebuie deplasată la dreapta cu 2^i biți. Se denumește mantisă inferioară, mantisa asociată exponentului mai mic. Cantitatea totală de deplasare a mantisei inferioare este deci dată de rezultatul diferenței celor doi exponenți. Totuși biții de ordin 5, 6 și 7 de valoare 1 ar implica deplasări de 32, 64 și respectiv 128 de biți ceea ce este imposibil întrucât mantisa are doar 23 de biți în simplă precizie. Aceasta înseamnă că dacă cel puțin unul din acești biți este 1, toți biții mantisei inferioare vor deveni egali cu zero.

Deoarece însă inițial nu se cunoaște care este semnul diferenței exponenților, nu se poate stabili care este mantisa inferioară, pentru a se putea starta deplasarea acesteia. Circuitul care realizează deplasarea mantisei, “barrel shifter”, introduce zerouri în stânga mantisei, pe măsură ce aceasta se deplasează spre dreapta. El este compus din mai multe etaje, sau selectoare, succesive. Selectorul de rang i fie lasă data nedeplasată, fie o deplasează cu 2^i poziții, în funcție de valoarea bitului s_i din diferența pozitivă a exponenților. O secțiune din blocul “barrel shifter” este prezentată în figura 3.1.

Pentru a evita reconversia din cod complement de doi în cod semn-mărime în cazul diferenței negative a celor doi exponenți $ExpA$ și $ExpB$ care sunt ambii pozitivi (având inclus deplasamentul de 127), se calculează simultan $ExpA-ExpB$ și $ExpB-ExpA$ și se reține diferența pozitivă.

De asemenea selectoarele din blocul “barrel shifter” se aranjează astfel încât să înceapă să opereze cu deplasări mici spre deplasări mari, ținând cont că și biții diferenței exponenților se stabilesc dinspre LSB spre MSB. Soluția pe care am propus-o ca și **varianta 1**, similară cu cea prezentată în [91], a fost utilizarea a două circuite “barrel shifter” care operează în paralel și care sunt comandate de biții diferențelor $ExpA-ExpB$, respectiv $ExpB-ExpA$. Astfel cele două mantise, a căror valoare inițială este conservată, încep să se deplaseze simultan cu stabilirea biților celor două diferențe. În final doar una din cele două date deplasate va fi utilă. Ea va fi selectată împreună cu cealaltă mantisă nedeplasată, în momentul în care este cunoscut rezultatul pozitiv al diferenței celor doi exponenți. Schema bloc a circuitului de aliniere a mantiselor este prezentată în figura 3.2.

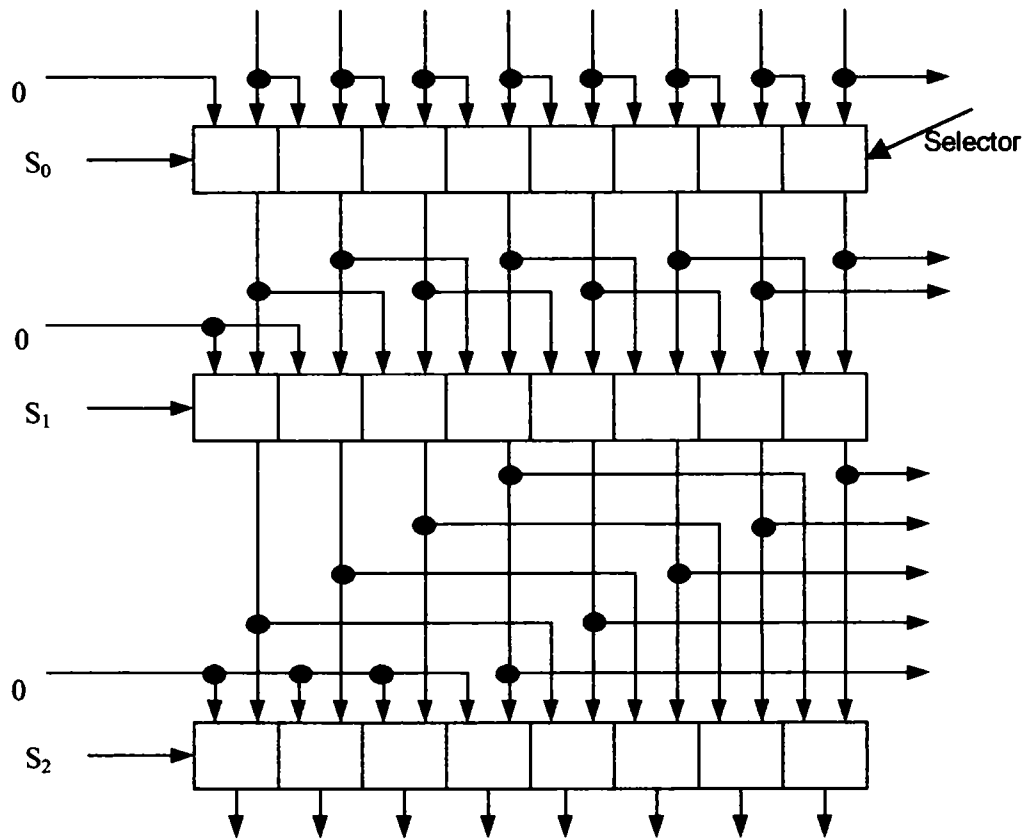


Fig. 3.1 Secțiune din blocul "barrel shifter"

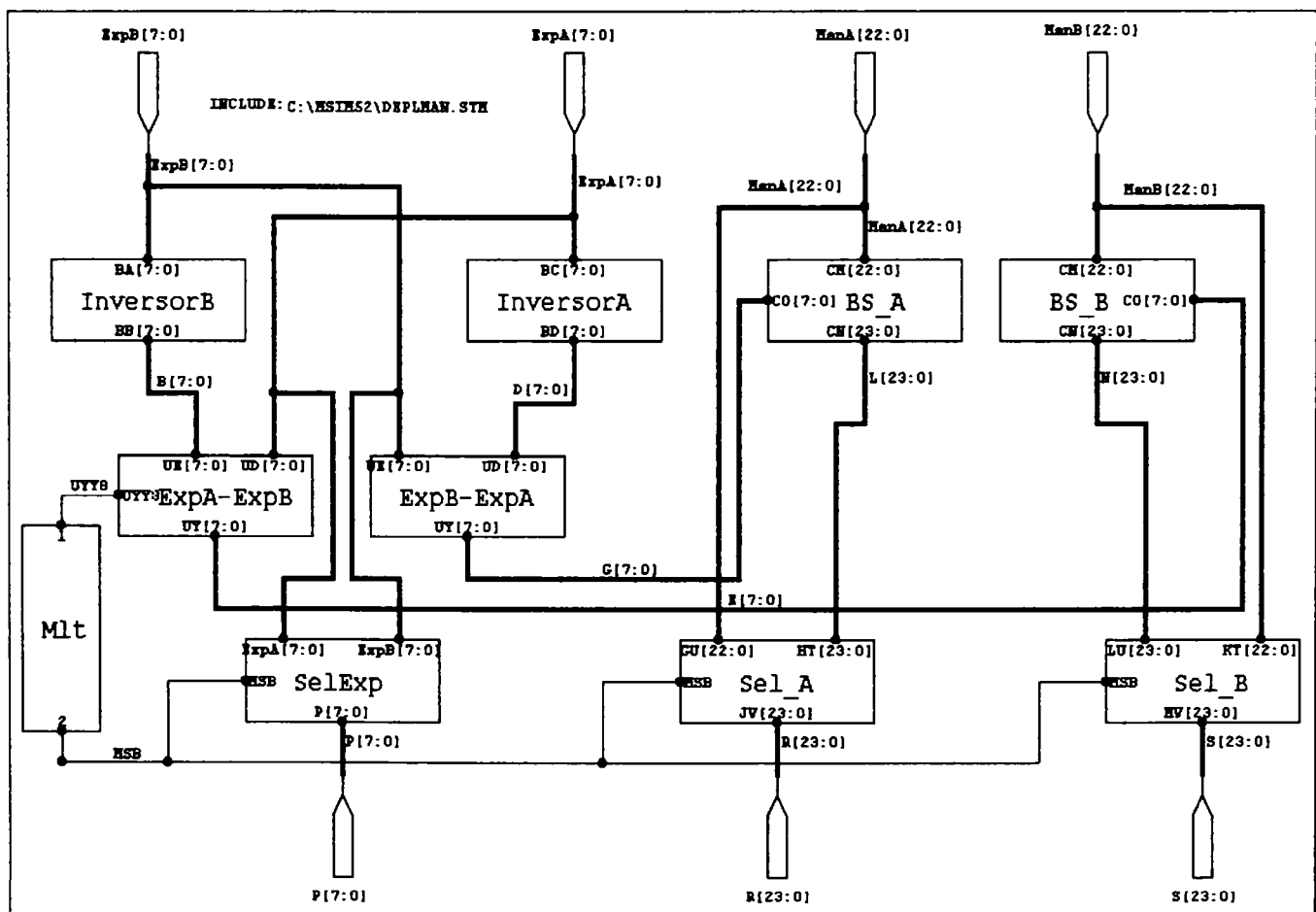


Fig.3.2 Schema bloc a circuitului de aliniere a mantiselor – varianta 1

Diferențele $ExpA-ExpB$, respectiv $ExpB-ExpA$ sunt efectuate folosind reprezentarea în complement de doi, cu ajutorul blocurilor InversorA, InversorB și al sumatoarelor “ $ExpA-ExpB$ ”, respectiv “ $ExpB-ExpA$ ” prin intermediul cărora se calculează $ExpA+1/ExpB+1$, respectiv $ExpB+1/ExpA+1$ ($/ExpA(B)$ este data obținută prin inversarea biților lui $ExpA(B)$). Sumatorul “ $ExpA-ExpB$ ” a fost extins la 9 biți pentru ca bitul MSB (UY8 în figura 3.2) al sumei obținute să poată sesiza semnul diferenței celor 2 exponenți pozitivi, extinși și ei la 9 biți prin adăugarea la stânga a câte unui zero. Spre deosebire de soluția prezentată în [91], cele două sumatoare folosite sunt sumatoare de concepție proprie, organizate pe două niveluri ierarhice, similare cu sumatorul din figura 2.12. În figura 3.3 este prezentată schema bloc a sumatorului rapid “ $ExpA-ExpB$ ”. Nivelul ierarhic inferior este reprezentat de trei sumatoare CLA (“carry look-ahead”) pe 4 biți (CLA1 cu transport inițial 1, necesar pentru efectuarea diferenței exponenților) și respectiv 5 biți (CLA2 cu transport inițial 1 și CLA3 cu transport inițial 0). Nivelul ierarhic superior implementează un mecanism de selecție a transportului prin care blocul SEL selectează rezultatul de la CLA2 sau CLA3 în funcție de transportul $c4$ generat de CLA1. Blocurile A', B' și B'' au fost prezentate în figura 2.10, iar blocul CON este conector de magistrală.

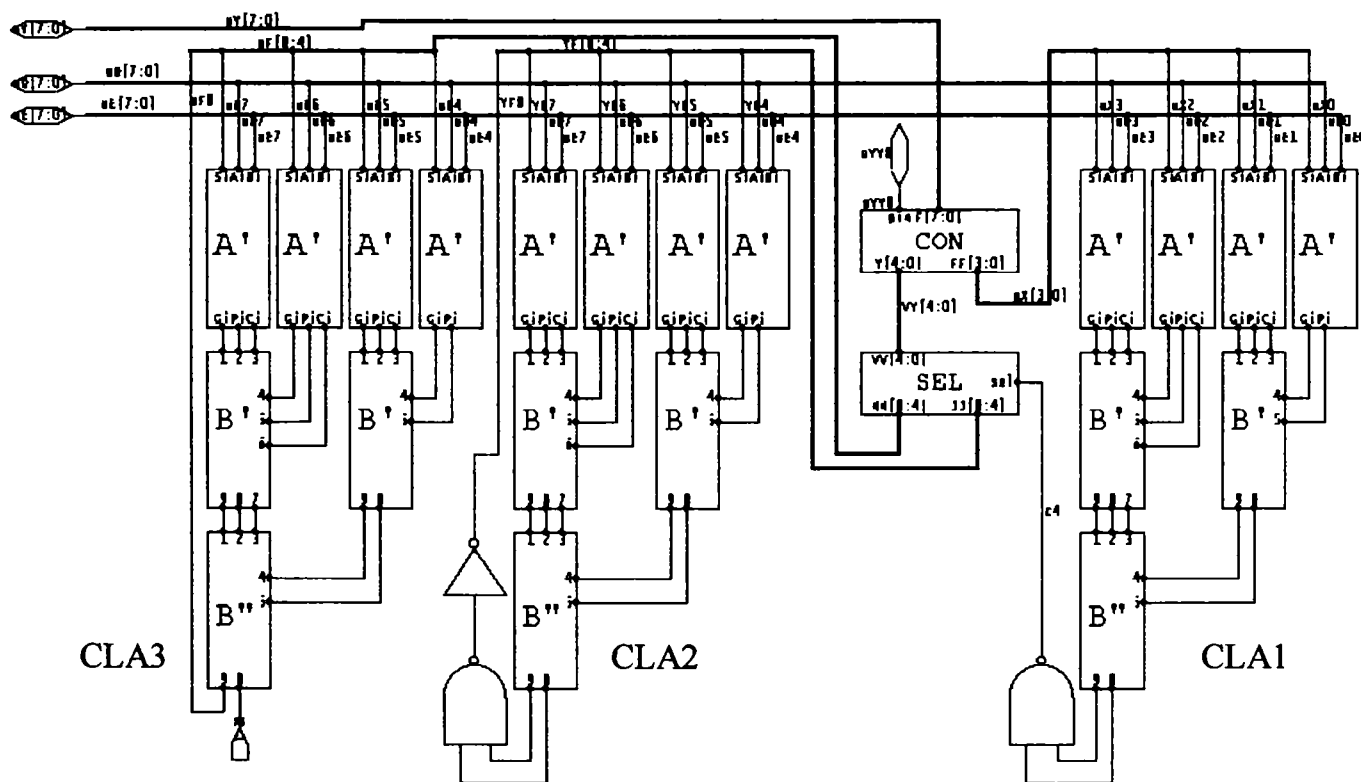


Fig. 3.3 Schema bloc a sumatorului rapid “ $ExpA-ExpB$ ”

Cele două blocuri “barrel shifter” din figura 3.2, BS_A și BS_B, conțin câte 6 selectoare de genul celor din figura 3.1, dar cu 24 de celule fiecare. Primul selector din fiecare circuit “barrel shifter” va avea la intrare obligatoriu MSB=1, adică partea întreagă a mantisei ce este astfel restituită părții fracționare pe 23 de biți. Sunt necesare doar 6 selectoare și nu 8 deoarece, așa cum s-a precizat anterior,

deplasări de 32, 64, respectiv 128 biți nu au sens, astfel că biții s_5 , s_6 și s_7 ai diferențelor exponenților au fost grupați printr-un circuit SAU pe o singură linie, notată s_{5+} . Selectoarele au fost proiectate utilizând porți de transmisie, respectiv tranzistoare de trecere (“pass-transistors”) complementare, pentru selectarea uneia sau alteia din cele două căi de acces ale fiecărei celule din selector. În acest fel fiecare linie de comandă, de la s_0 la s_{5+} ar trebui să suporte încărcarea capacitivă a $2 \times 24 = 48$ grile de tranzistoare MOS. Acest lucru ar fi dezastruos atât pentru porțile logice ce furnizează biții s_i cât și pentru viteza de operare. Problema se pune și mai pregnant pentru linia MSB din figura 2.2 care trebuie să comande 112 grile de tranzistoare MOS.

Evitarea acestor situații s-a făcut utilizând structuri arborescente de inversoare CMOS, pentru multiplicarea capabilității de comandă a liniilor menționate. Arborele de la capătul fiecărei linii s_i va conține patru etaje de inversoare ce nu introduc însă o întârziere semnificativă.

În cazul liniei MSB, în blocul “Mlt” din figura 2.2, se va utiliza un arbore cu 5 niveluri, la care ultimul nivel va avea 32 ieșiri CMOS, care vor fi folosite pentru comanda celor 96 grile MOS din Sel_A și Sel_B și pentru cele 16 grile din Sel_Exp. Astfel fiecărei ieșiri CMOS i s-au considerat alocate maxim 4 grile de tranzistoare.

Bitul MSB de la “ExpA-ExpB” este cel care selectează la ieșirea etajului exponentul mai mare cu ajutorul blocului SelExp, respectiv perechea corectă mantisă nedepasată - mantisă depasată, cu ajutorul blocurilor Sel_A și Sel_B. Logica de selecție este următoarea: dacă MSB este 1, atunci va fi selectat ExpB și rezultatul furnizat de blocul BS_A, iar dacă MSB este 0, atunci va fi selectat ExpA și rezultatul furnizat de BS_B.

Toate blocurile componente și subblocurile din schemele ierarhice inferioare au fost proiectate până la nivel de poartă logică sau poartă de transmisie și sunt prezentate în Anexa 5.

Au fost efectuate simulări pentru cazul cel mai defavorabil și a fost obținut un timp maxim de propagare prin circuitul din figura 2.2 de 2,4 ns (utilizând etalonul - tehnologie CMOS 0,5 microni). Acest rezultat recomandă circuitul proiectat pentru procesoarele aritmetice în virgulă flotantă, însă nu e potrivit cu ceea ce se dorește în cazul de față, adică timp de propagare aproape de 4ns și arie ocupată de circuit cât mai mică.

Într-adevăr aria ocupată pe cip este mare deoarece două circuite “barrel shifter” înseamnă nu numai număr dublu de tranzistoare ci și spațiu de rutare între selectoare dublu. Într-adevăr, așa cum se poate observa chiar și din figura 2.1, spațiul de rutare pentru deplasare cu 2^i este proporțional cu 2^i , adică devine din ce în ce mai mare pe măsură ce se accede spre selectoarele care realizează deplasări mari. La acest spațiu se adaugă și cel necesar pentru multiplicarea capabilității de comandă a liniei de selecție s_i , același pentru fiecare selector. De asemenea cele două sumatoare rapide, precum și cele trei blocuri selectoare ocupă o suprafață importantă.

Pentru a realiza o mai bună apropiere de scopul propus, la elaborarea variantei 2 de implementare a circuitului de aliniere a mantiselor, am plecat de la soluția descrisă în [26] la nivel de schemă bloc și pe care am implementat-o (cu anumite modificări) printr-o proiectare completă până la nivel de poartă logică. Rezultatele obținute au fost prezentate în [39]. În această variantă se utilizează trei circuite “barrel shifter”, două care operează deplasări mici pentru fiecare mantisă, până la maxim 7 poziții (conțin câte 3 selectoare) și un al treilea care primește deja mantisa corectă, parțial deplasată și care conține tot trei selectoare. În acest fel, spațiul de rutare excesiv, necesar pentru realizarea deplasărilor mari s-a redus la jumătate. Schema bloc a acestui circuit este prezentată în figura 3.4.

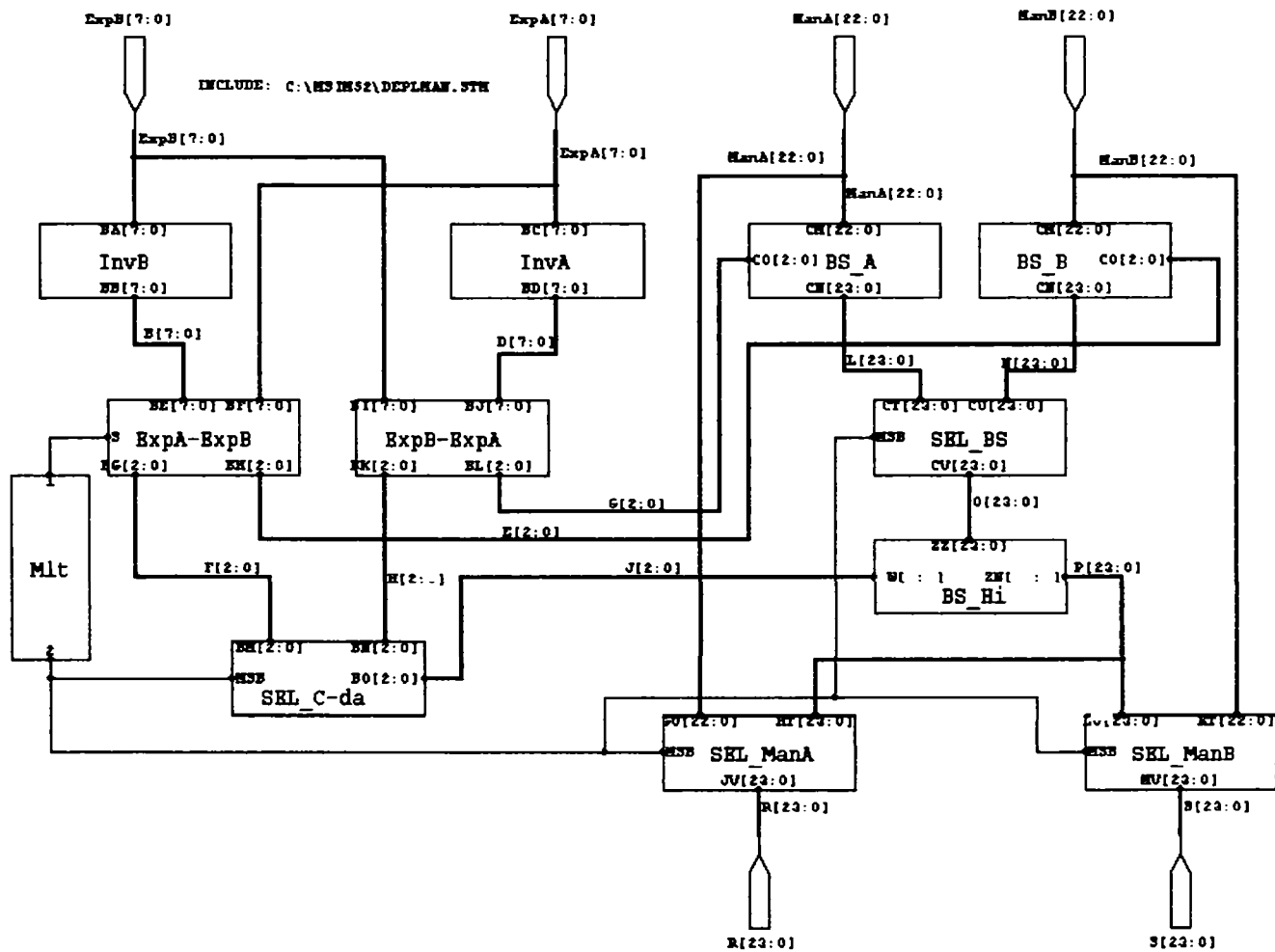


Fig. 3.4 Schema bloc a circuitului de aliniere a mantiselor – varianta 2

Pentru simplitate, din schema bloc lipsește circuitul de selecție a exponentului mai mare, dar acest lucru nu afectează viteza de operare deoarece circuitul menționat nu se găsește pe calea critică de propagare.

Cele 6 linii de comandă (de la s_0 la s_5) ale circuitelor “barrel shifter”, furnizate de blocurile “ExpA-ExpB” și “ExpB-ExpA”, sunt grupate în două perechi de câte 3 biți. Perechea ce conține cei trei biți mai puțin semnificativi ($s_0 \div s_2$) de la ieșirea fiecărui sumator, este folosită pentru comanda circuitelor BS_A și BS_B care realizează deplasările mici. Din cealaltă pereche de triplete de biți se va selecta în blocul “Sel_C-dă” cu ajutorul aceluiași MSB de la sumatorul “ExpA-ExpB”, tripleta care va comanda circuitul “barrel shifter” BS_Hi. În același

În această variantă este folosit un singur circuit “barrel shifter”, (blocul BS din figura 3.5) care deplasează mantisa inferioară, stabilită după aflarea semnului diferenței celor doi exponenți. Nu poate fi folosit însă un sumator unic pentru calculul unei singure diferențe a exponenților întrucât, dacă rezultatul se dovedește a fi negativ el va fi reprezentat în complement de doi și trebuie reconvertit în cod mărime-semn, ceea ce întârzie prea mult producerea rezultatului.

Linia MSB de la ieșirea sumatorului “ExpA-ExpB”, identic cu cel folosit în varianta 2, realizează următoarele funcțiuni:

- selectează mantisa asociată exponentului mai mic, adică mantisa inferioară, pentru a fi aplicată la intrarea blocului BS;
- selectează prin intermediul blocului Sel_C-dă liniile $s_0 \div s_{5+}$ furnizate de sumatorul al cărui rezultat este pozitiv, pentru a comanda deplasările mantisei inferioare în blocul BS;
- selectează mantisa asociată exponentului mai mare, adică mantisa superioară, direct la ieșirea circuitului. Prin intermediul blocurilor Sel_ManA și Sel_ManB este păstrată poziția relativă a celor două mantise după alinierea lor;
- selectează prin intermediul blocului SelecExp exponentul mai mare la ieșirea circuitului.

Ca și în variantele anterioare, capacitatea de comandă a liniei MSB a trebuit să fie multiplicată folosind un arbore de inversoare CMOS, ceea ce a implicat introducerea în schema simulată a blocului Mlt (figura 3.5).

Un alt aspect important de care s-a ținut cont la construcția circuitului “barrel shifter” a fost limitarea capacităților parazite care se acumulează datorită cascaderii porților de transmisie din circuite selectoare consecutive și care încarcă excesiv latch-ul de date de la intrarea acestui nivel pipeline. Din acest motiv blocul BS a fost secționat de un rând dublu de inversoare CMOS, care nu modifică valoarea mantisei inferioare parțial deplasate ce se aplică mai departe următoarelor selectoare din circuitul “barrel shifter”. În acest fel fronturile semnalelor propagate prin structură sunt refăcute și nu apar întârzieri critice.

În urma simulărilor efectuate, s-a obținut un timp maxim de propagare pe acest nivel pipeline de 4 ns. Schemele detaliate ale circuitului sunt de asemenea prezentate în Anexa 5. Economia de suprafață pe cip a fost estimată la aproape 50% față de varianta 1.

3.2 Proiectarea circuitului de adunare/scădere a mantiselor aliniat

Circuitul de adunare/scădere a mantiselor aliniat operează pe 25 de biți întrucât la cei 23 de biți memorați ai mantiselor s-a restituit bitul implicit din stânga virgulei, precum și bitul de semn. *Circuitul este de același tip cu cel utilizat în unitatea logaritmică și prezentat în figura 2.9, deci implementează în același mod operația de scădere, dar este mai complex întrucât include blocuri inversoare pentru ambii operanzi (oricare din ei poate fi negativ) și un bloc logic care controlează selectoarele din blocurile inversoare, respectiv selectorul final, în*

funcție de operația ce trebuie efectuată (adunare sau scădere) cât și de semnul celor doi operanzi. Schema lui bloc este prezentată în figura 3.6.

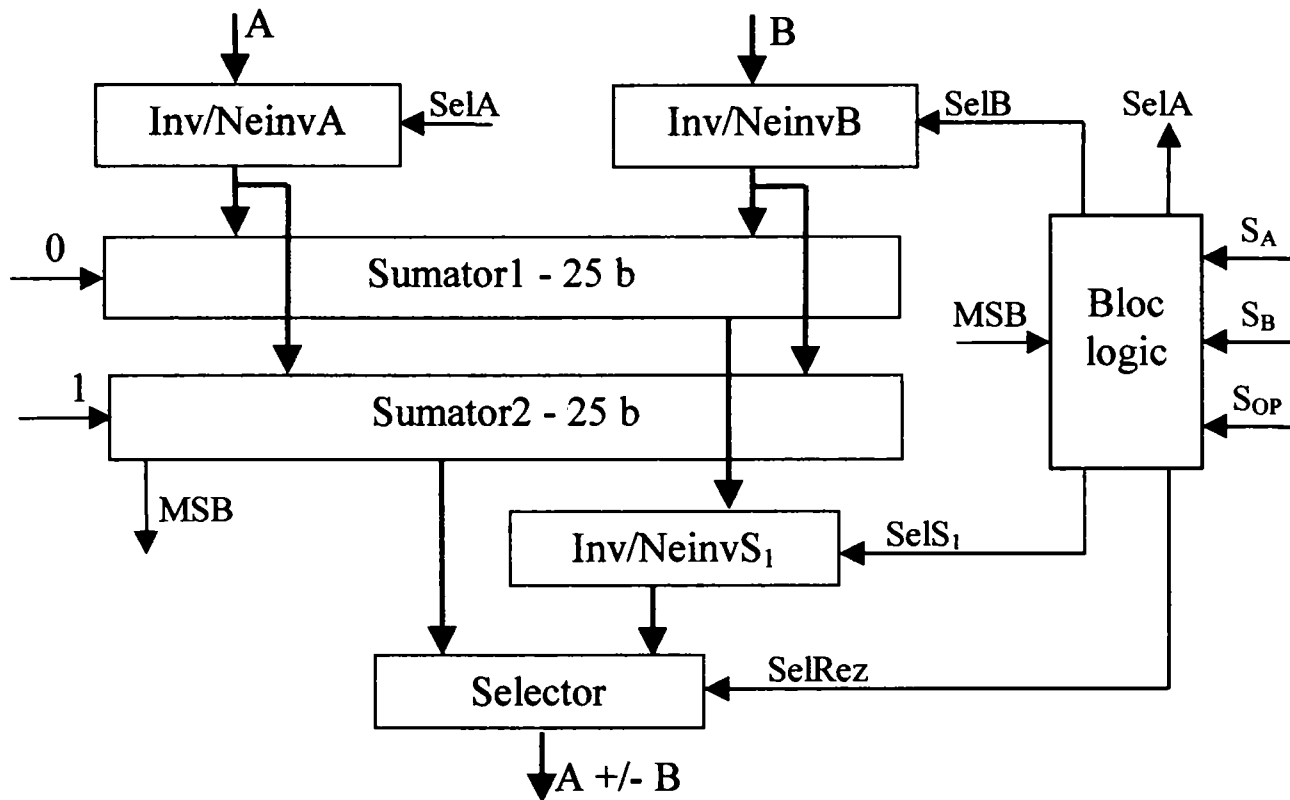


Fig. 3.6 Schema bloc a circuitului de adunare/scădere a mantiselor aliniat

Considerând A și B ca fiind cele două mantise aliniat, S_A și S_B semnele lor, iar S_{OP} bitul asociat operației (0 pentru adunare și 1 pentru scădere), atunci blocul logic implementează următoarea regulă:

-dacă:

$$(a1): S_A = 0, S_B = 0, S_{OP} = 0, \text{ sau}$$

$$(a2): S_A = 1, S_B = 1, S_{OP} = 0, \text{ sau}$$

$$(a3): S_A = 0, S_B = 1, S_{OP} = 1, \text{ sau}$$

$$(a4): S_A = 1, S_B = 0, S_{OP} = 1,$$

atunci se efectuează adunare propriu-zisă, iar semnul rezultatului este semnul lui A ($S_{Rez} = S_A$). În acest caz: $S_A \oplus S_B \oplus S_{OP} = 0$.

-dacă:

$$(s1): S_A = 0, S_B = 0, S_{OP} = 1, \text{ sau}$$

$$(s2): S_A = 1, S_B = 1, S_{OP} = 1, \text{ sau}$$

$$(s3): S_A = 0, S_B = 1, S_{OP} = 0, \text{ sau}$$

$$(s4): S_A = 1, S_B = 0, S_{OP} = 0,$$

atunci se efectuează scădere propriu-zisă pe 25 de biți. Se constată că în acest caz: $S_A \oplus S_B \oplus S_{OP} = 1$.

În cazul adunării propriu-zise, blocurile $Inv/NeinvA$, $Inv/NeinvB$ și $Inv/NeinvS_1$ trebuie să fie transparente ($SelA = SelB = SelS_1 = 0$) iar blocul Selector va selecta rezultatul furnizat de Sumator1 ($SelRez = 0$).

În cazul scăderii propriu-zise, blocul Inv/NeinvA va fi transparent în situațiile (s1) și (s3) pentru care linia SelA=0, iar blocul Inv/NeinvB va fi transparent în situațiile (s2) și (s4) pentru care linia SelB=0. Blocul Inv/NeinvS₁ va efectua inversarea biților de la Sumator1 (SelS₁=1), iar blocul Selector va selecta prin semnalul SelRez rezultatul de la Sumator2 dacă MSB=1, respectiv de la blocul Inv/NeinvS₁ dacă MSB=0. Semnul rezultatului este în acest caz dat de MSB. În concluzie Blocul logic trebuie să asigure, în concordanță cu constatările anterioare, următoarele funcții:

$$SelA = S_A S_B S_{OP} + S_A /S_B /S_{OP} = S_A (S_B \oplus S_{OP}) \quad (3.1)$$

$$SelB = /S_A /S_B S_{OP} + /S_A S_B /S_{OP} = /S_A (S_B \oplus S_{OP}) \quad (3.2)$$

în care /S_A, /S_B și /S_{OP} sunt S_A negat, S_B negat, respectiv S_{OP} negat. De asemenea Blocul logic mai furnizează:

$$SelS_1 = S_A \oplus S_B \oplus S_{OP} \quad (3.3)$$

$$SelRez = /MSB (S_A \oplus S_B \oplus S_{OP}) \quad (3.4)$$

Schema blocului logic este prezentată în figura 3.7. Așa cum se observă din această figură, blocul logic a fost astfel proiectat încât nici o ieșire să nu comande mai mult de două intrări logice, respectiv mai mult de 4 grile MOS. Blocul logic mai furnizează prin intermediul ieșirii Srez, semnul rezultatului final. Pe ansamblul circuitului de adunare/scădere, calea critică de propagare trece și prin interiorul acestui bloc și cuprinde linia MSB. Din acest motiv, căile care prelucrează bitul MSB includ un număr minim de niveluri logice.

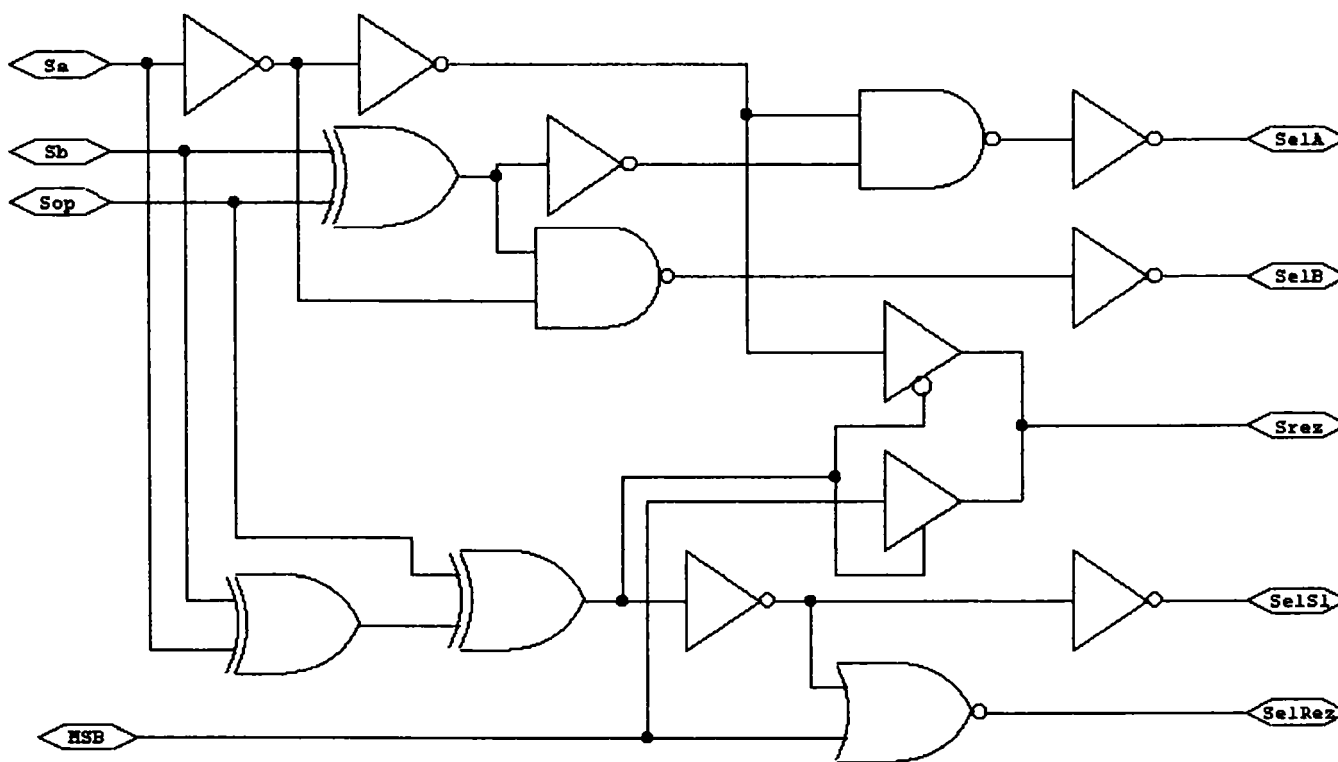


Fig. 3.7 Schema detaliată a Blocului logic

În *varianta 1* de realizare a circuitului de adunare/acădere blocurile Sumator1 și Sumator2 din figura 3.6, cu transport inițial 0 respectiv 1, sunt sumatoare cu selecția transportului, la care secțiunile componente de câte 8 biți sunt realizate cu sumatoare de tip CLA. Aceste blocuri au structura similară cu cea a sumatorului prezentat în figura 2.12, beneficiind astfel de toate îmbunătățirile aduse acestuia în privința vitezei de operare. Schema bloc a acestui tip de sumator, de concepție proprie este prezentată în figura 3.8.

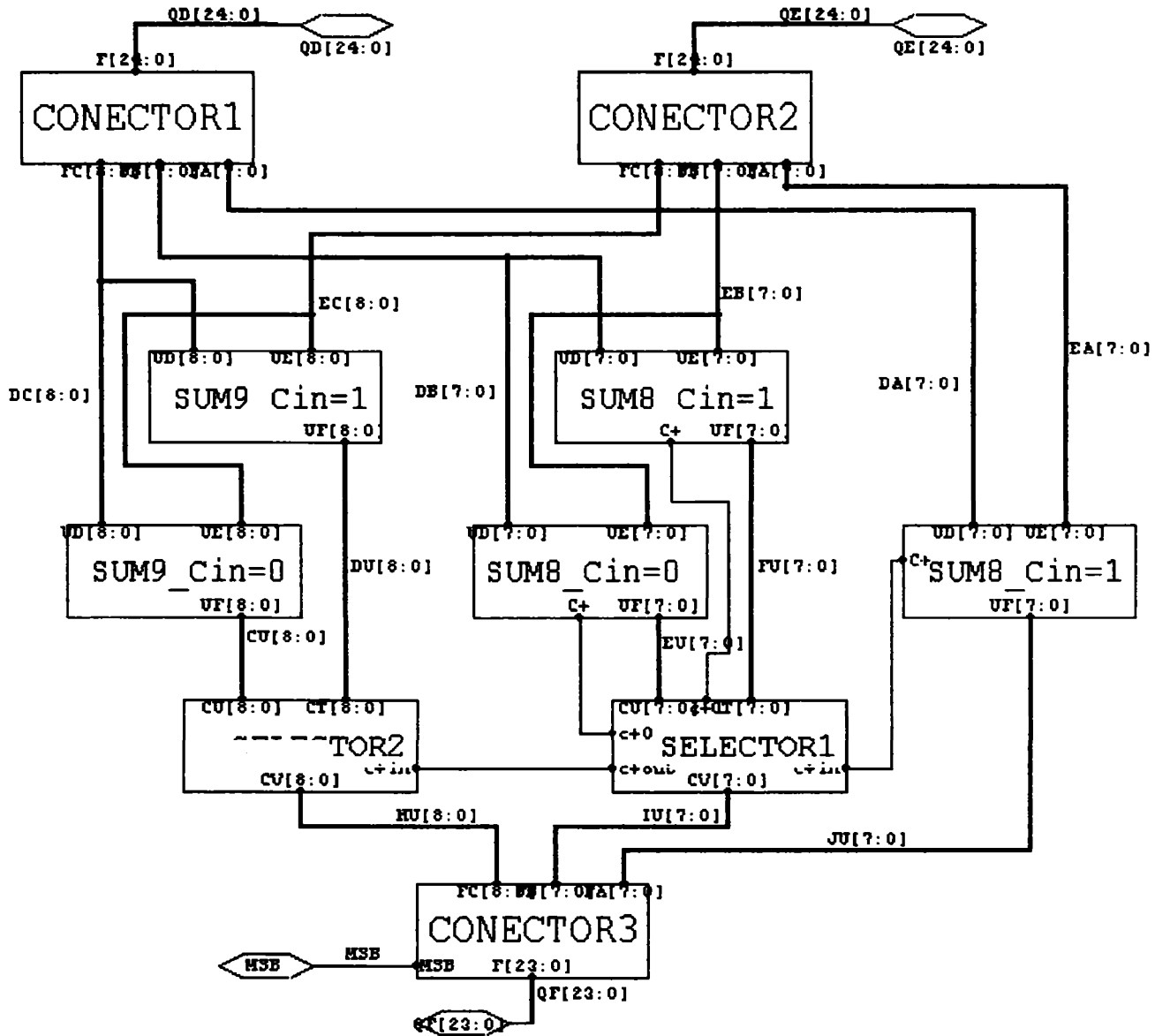


Fig. 3.8 Schema bloc a sumatorului pe 25 de biți – varianta 1

Bitul al 25-lea (MSB) al sumei, în concordanță cu observația făcută în paragraful 2.2.2.4 la figura 2.11, se generează mai repede decât al 24-lea bit și deci nu conduce la creșterea timpului de operare al ultimei secțiuni a sumatorului, care va fi pe 9 biți.

În figura 3.6 blocurile Inv/NeinvA, Inv/NeinvB, Inv/NeinvS₁ și Selector includ arbori de inversoare CMOS pentru multiplicarea capabilității de comandă a liniilor de selecție SelA, SelB, SelS₁ și SelRez. Multiplicarea capabilității de comandă a liniilor SelA, SelB și SelS₁ nu afectează viteza sumatorului deoarece S_A,

S_B și S_{OP} determină starea lor logică încă din timpul operației de aliniere a mantiselor. În schimb, în cazul liniei SelRez, care depinde și de MSB, viteza sumatorului este întrucâtva afectată. Totuși acest efect este compensat parțial de faptul că în timpul multiplicării capabilității de comandă a liniei SelRez se efectuează inversarea și selectarea inversării biților de la Sumator1 în blocul Inv/NeinvS₁.

Arhitectura propusă pentru circuitul de adunare/scădere prezentată în figura 3.6, diferită de cea prezentată în [26], oferă următoarele avantaje față de aceasta:

- datorită inversării ordinii celor două blocuri finale, Inv/NeinvS₁ și Selector, se elimină timpul necesar inversării și selecției inversării pentru rezultatul final;
- prezența Blocului logic în interiorul circuitului de adunare/scădere permite adunarea și scăderea operanzilor în toate cazurile posibile din punct de vedere al semnului lor și nu doar adunare de operanzi de același semn sau de semne contrare, cum se procedează în anumite unități de calcul;
- controlul circuitului de adunare/scădere este foarte eficient și este astfel conceput încât pentru generarea semnalelor necesare numărul nivelurilor logice parcurse este minim; în plus Blocul logic poate fi distribuit între primul și al doilea nivel pipeline al subunității de calcul în virgulă flotantă, ceea ce permite scurtarea semnificativă a căii critice de propagare;
- utilizarea sumatoarelor cu selecția transportului, la care secțiunile componente de câte 8 biți sunt realizate cu sumatoare de tip CLA îmbunătățite, a condus la o creștere substanțială a vitezei de calcul.

În urma simulărilor efectuate, considerând tehnologie CMOS 0,5 μm, timpul de propagare pentru cazul cel mai defavorabil a fost găsit de 4 ns. În realitate însă, timpul total de calcul pentru circuitul de adunare/scădere este mult mai scurt deoarece valoarea logică a liniilor SelA și SelB din figura 3.6 poate fi cunoscută încă din timpul etapei de deplasare a mantiselor. De asemenea, propagarea transportului prin arborele de multiplicare a capabilității de comandă a acestor linii poate fi realizată în etapa precedentă. În concluzie, timpul total de operare este de doar 2,9 ns. Acest rezultat recomandă utilizarea circuitului la concepția unor unități de calcul performante în care se operează exclusiv în virgulă flotantă.

În varianta 2 de realizare a circuitului de adunare/scădere se păstrează aceeași arhitectură avantajoasă și se exploatează eficiența Blocului logic, de această dată nu pentru a obține o viteză maximă de calcul, ci pentru a putea utiliza, pentru adunarea/scăderea mantiselor aliniolate, sumatoare care necesită o suprafață mai mică pe cip. Condiția care a trebuit respectată a fost ca timpul de producere al rezultatului să se încadreze în 4 ns. Prin utilizarea unui sumator CLA îmbunătățit, de tipul celui din figura 2.12, dar pe 25 de biți, s-a confirmat un timp de propagare pentru cazul cel mai defavorabil de 4 ns. Numărul porților logice și al porților de transmisie a scăzut de la circa 700 la aproximativ 350 ceea ce înseamnă că s-a obținut 50% economie de suprafață față de varianta 1. Schemele de detaliu ale celor două variante de circuit sumator/scăzător se găsesc în Anexa 6.

3.3 Circuitul de normalizare

Circuitul de normalizare ocupă al treilea nivel pipeline al subunității de adunare/scădere în virgulă flotantă și are schema bloc prezentată în figura 3.9.

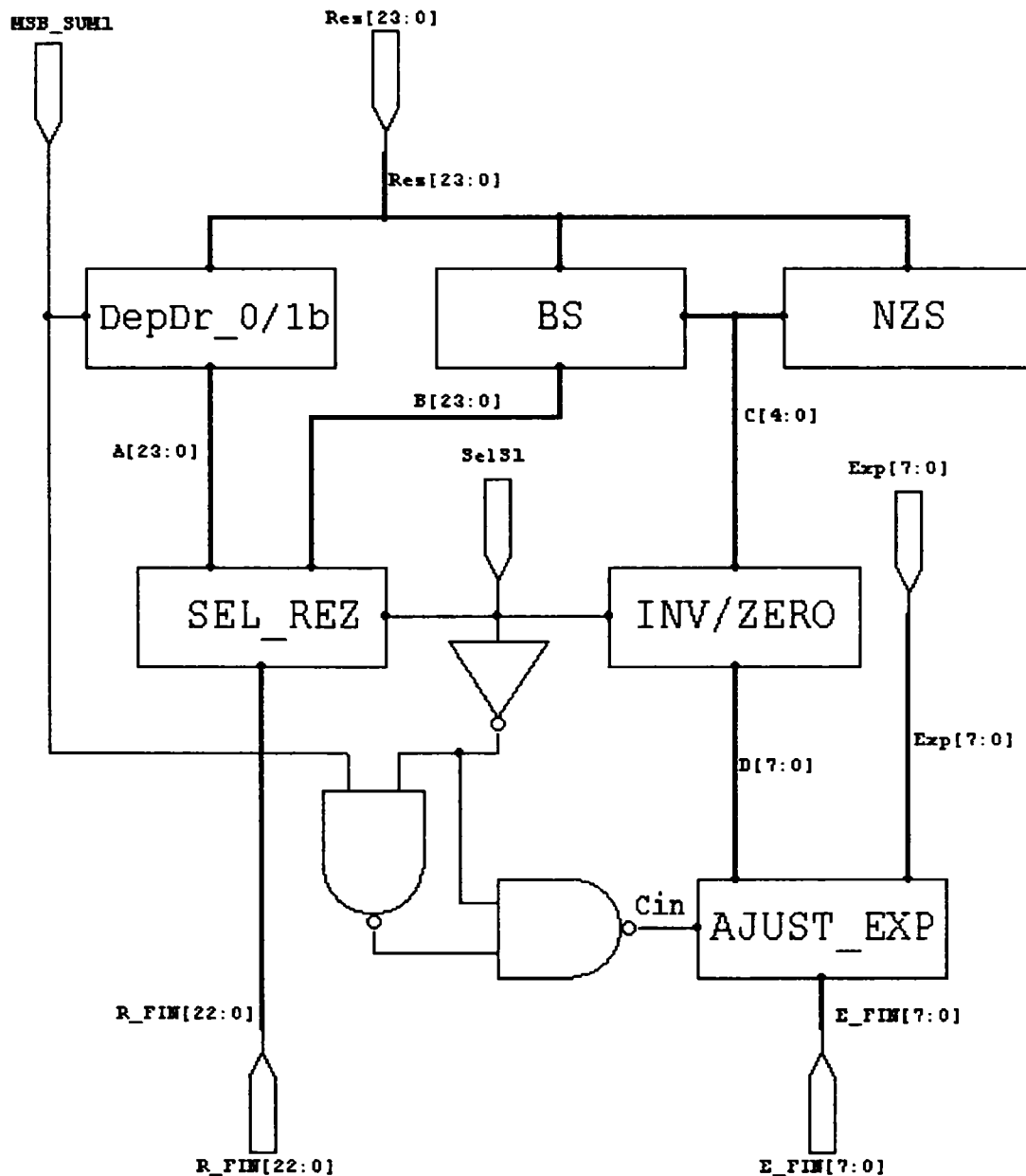


Fig. 3.9 Schema bloc a circuitului de normalizare

În urma efectuării unei adunări propriu-zise, rezultatul în valoare absolută se poate găsi în intervalul $[0,4)$ ceea ce înseamnă că operația de normalizare e foarte simplă. Rezultatul fie se lasă nedeplasat, fie se deplasează la dreapta cu un bit, prin intermediul blocului “DepDr_0/1b”, bineînțeles în funcție de valoarea bitului MSB_SUM1 de la blocul Sumator1 din figura 3.6. Atunci însă când se efectuează o scădere propriu-zisă, dacă operanzii sunt foarte apropiați în valoare absolută, pot rezulta multe zerouri înaintea celui mai semnificativ 1 din rezultat. Aceste zerouri se denumesc zerouri semnificative. Pentru a se efectua normalizarea, zerourile semnificative sunt în primă fază contorizate în blocul NZS, urmând ca apoi să fie comandat circuitul “barrel shifter” BS, similar cu cel utilizat la deplasarea mantisei

inferioare, în vederea normalizării. Evident în acest caz are loc o pierdere inevitabilă de precizie, întrucât deplasările succesive la stânga ale rezultatului diferență implică introducerea de zerouri dinspre dreapta în pozițiile rămase libere.

În final, blocul SEL_REZ va selecta rezultatul de la ieșirea blocului "DepDr_0/1b" sau de la ieșirea blocului BS în funcție de faptul că în nivelul pipeline 2 s-a efectuat o adunare propriu-zisă (SelS1 din figura 3.6 și dat de relația 3.3 - egal cu zero) sau o scădere propriu-zisă (SelS1=1).

În același timp, exponentul transmis de la ieșirea primului nivel pipeline al subunității de adunare/scădere în virgulă flotantă, se ajustează în nivelul 3 în funcție de rezultatul normalizării mantisei furnizate la ieșirea nivelului 2. Pentru a realiza acest lucru, în mod normal sunt necesare două circuite sumatoare pe 9 biți. Unul asigură incrementarea exponentului cu o unitate, în cazul unei adunări propriu-zise în care rezultă doi biți în stânga virgulei (MSB_SUM1=1). Celălalt asigură scăderea a n unități din exponent, în cazul unei scăderi propriu-zise în care rezultă n zerouri semnificative în mantisa furnizată de nivelul pipeline 2. În final este selectat rezultatul corespunzător operației efectuate.

O contribuție proprie a fost adusă și în cazul acestui circuit, concretizată prin utilizarea unui singur circuit sumator pe 9 biți, blocul AJUST_EXP, așa cum se observă în figura 3.9. Pe baza semnalelor MSB_SUM1 și SelS1 furnizat de Blocul logic din nivelul 2 (figura 3.7) se asigură un transport inițial pentru sumatorul AJUST_EXP în următoarele două situații:

- atunci când se efectuează scădere propriu-zisă și trebuie implementată relația:

$$\text{Exp} - \text{Nzs} = \text{Exp} + \text{/Nzs} + 1 \quad (3.5)$$

în care Exp este exponentul furnizat de nivelul 1 al subunității în virgulă flotantă și la care se atașează un zero ca bit de semn, Nzs este numărul zerourilor semnificative exprimat în cod mărime-semn pe 9 biți, iar /Nzs este data obținută prin inversarea biților lui Nzs cu ajutorul blocului INV/ZERO (figura 3.9) atunci când SelS1=1; dacă rezultatul obținut la ieșirea acestui bloc rezultă negativ (bitul $\text{Exp}_8=1$), în ciuda deplasamentului de 127 inclus, se semnalează subdepășire;

- atunci când se efectuează adunare propriu-zisă și trebuie incrementat cu o unitate exponentul, întrucât prin însumarea celor două mantise au rezultat doi biți în stânga virgulei și mantisa a fost deplasată spre dreapta cu un bit; în acest caz (SelS1=0) la ieșirea blocului INV/ZERO este furnizat un cuvânt nul și se va calcula suma $\text{Exp}+1$; dacă în acest caz rezultă bitul $\text{Exp}_8=1$, fiind vorba de adunarea a două numere pozitive, se semnalează supradepășire.

Astfel, pentru a obține semnalul Cin, ce reprezintă transportul inițial al sumatorului pentru ajustarea exponentului, a fost implementată relația:

$$\text{Cin} = \text{SelS1} + (\text{/SelS1})\text{MSB_SUM1} \quad (3.6)$$

În continuare se va face o analiză a timpului de propagare pe calea critică de propagare prin structura din figura 3.9, care include blocurile NZS, INV/ZERO și AJUST_EXP. Analiza se face în scopul determinării tipului de sumator ce poate

fi folosit pentru ca timpul de generare a exponentului ajustat să se încadreze în 4 ns, iar suprafața ocupată pe cip să fie minimă.

Blocul NZS a fost proiectat în detaliu după metoda descrisă în [85]. Schema sa bloc este prezentată în figura 3.10.

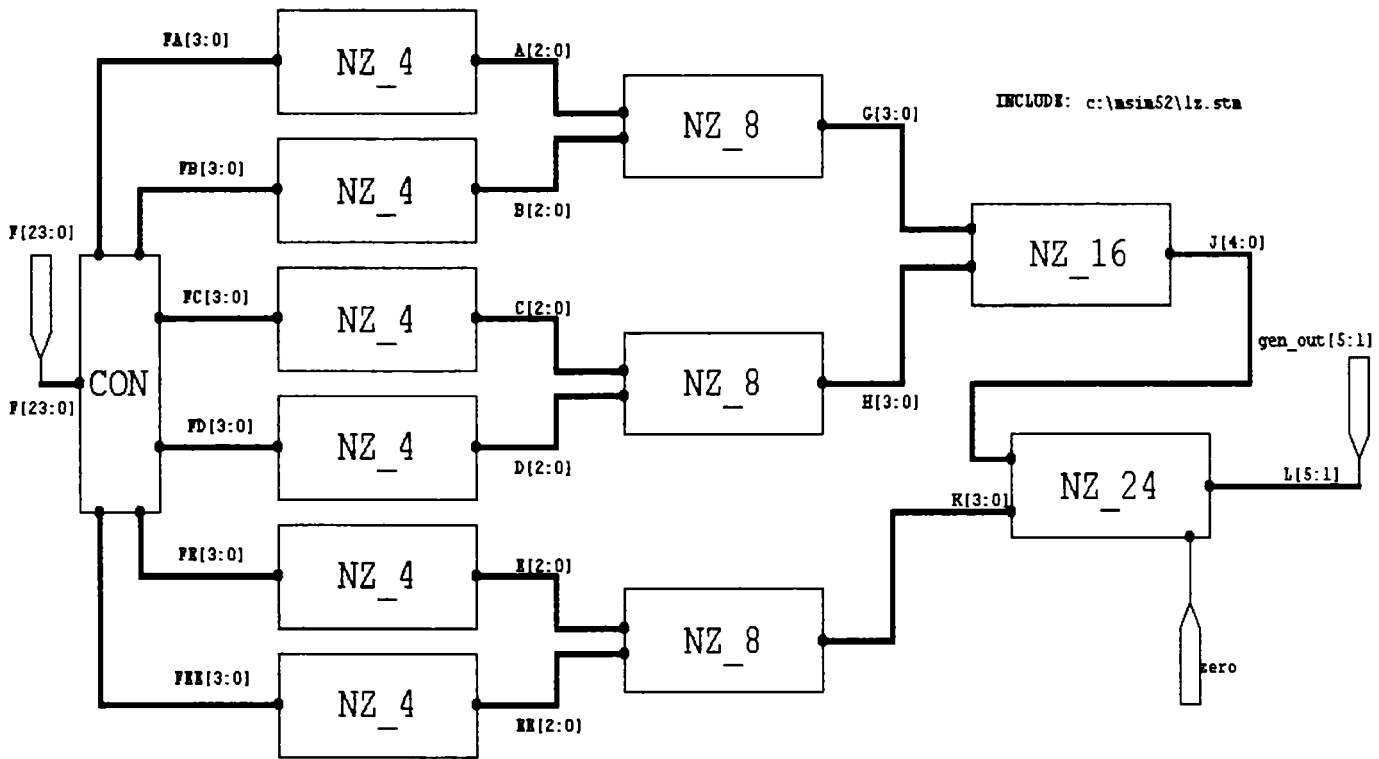


Fig. 3.10 Schema bloc a număratorului de zerouri semnificative (NZS)

Așa cum se observă în figura 3.10, circuitul NZS are o structură arborescentă, formată din patru tipuri de subblocuri, fiecare dintre acestea putând detecta un număr de zerouri semnificative (marcat în tipul subblocului) corespunzător grupelor de biți de intrare asociate lor. Blocul CON servește doar la disocierea magistralei de intrare de 24 de biți în 6 magistrale de câte 4 biți, prezența sa fiind impusă de simulator. Rangul celor 24 de biți crește de sus în jos. În fiecare subbloc există o poartă SAU și unul (în NZ_4), doi (în NZ_8), trei (în NZ_16), respectiv patru multiplexori cu 2 intrări (în NZ_24), prin care este implementat un mecanism prioritar de selecție. Poarta SAU din subblocul curent comandă multiplexorii subblocului din aval. La ieșirea subblocului NZ_24 se obține numărul reprezentat pe 5 biți, al zerourilor semnificative ale cuvântului de 24 de biți aplicat la intrarea circuitului NZS. Bitul "zero" furnizat de același subbloc indică situația în care toți biții de la intrare sunt nuli, caz în care operația de normalizare nu mai are sens, iar rezultatul adunării/scăderii celor 2 operanzi este declarat zero. Schemele detaliate ale tuturor blocurilor și subblocurilor componente se găsesc în Anexa 7.

Timpul critic de propagare printr-un subbloc este dat de poarta SAU (de fapt există porți SAU-NU alternate cu porți ȘI-NU, pentru a realiza funcții SAU cu porți intrinseci) și este de 0,2 ns în tehnologie CMOS de 0,5 μm . În concluzie, timpul de propagare prin întregul bloc este de 0,8 ns. La acest timp se mai adaugă

0,2 ns ale blocului INV/ZERO, astfel că pentru sumatorul de ajustare a exponentului rămân 3 ns din cele 4 ns alocate și acestui nivel pipeline. *Utilizarea unui sumator rapid pe 9 biți, identic cu cel prezentat în figura 3.3 ar conduce la obținerea unui timp de propagare pe acest nivel sub 3 ns. Utilizarea însă a unui sumator CLA modificat, de tipul celui din figura 2.11 dar pe 9 biți, asigură un timp de propagare de 2,5 ns în cazul cel mai defavorabil și reprezintă varianta optimă ce va fi adoptată în cazul de față. Timpul maxim de propagare prin ultimul nivel pipeline va fi deci de 3,5 ns.*

În concluzie, subunitatea de adunare/scădere în virgulă flotantă organizată pe trei niveluri pipeline este capabilă să funcționeze cu o perioadă de tact de 4 ns, ca și subunitatea logaritmică, în condițiile în care au fost utilizate circuite optimizate ce ocupă o suprafață minimă pe cip. Cele două subunități de calcul ale procesorului aritmetic permit astfel implementarea unei structuri de acumulator, pentru a se putea efectua rapid sume de produse de factori. În Anexa 8 este prezentat un exemplu de rulare al unui algoritm DSP simplu, prin care se efectuează o buclă Livermoore cu 20 de termeni (suma a 20 de produse de câte doi factori).

Cap.4 Rezultatele simulărilor blocurilor și subblocurilor componente ale procesorului hibrid în virgulă flotantă și logaritmic

În acest capitol vor fi prezentate rezultatele simulărilor majorității blocurilor și subblocurilor componente ale celor două subunități de calcul ale procesorului aritmetic hibrid conceput în lucrare. Au fost simulate astfel toate blocurile și subblocurile componente de pe căile critice de propagare, pentru a confirma contribuțiile aduse de autorul tezei și pentru a dovedi performanța sistemului de calcul proiectat. De asemenea au fost simulate și acele blocuri la care contribuțiile aduse s-au referit la optimizarea structurii lor din punct de vedere al suprafeței ocupate pe cip, atunci când aceste blocuri nu reprezentau o cale critică de propagare.

În cazul subunității logaritmice au fost făcute simulări pentru ambele variante prezentate. Au fost astfel evidențiate avantajele de viteză aduse prin implementarea metodei logaritmului mascat comparativ cu prima variantă și cu rezultatele obținute de F. S. Lai în [50], [51] și [107] ce au reprezentat ultimele realizări în domeniul procesoarelor hibride, prezentate de literatura de specialitate.

În cazul subunității de calcul în virgulă flotantă au fost de asemenea simulate toate variantele elaborate, pentru a putea fi depistată structura optimă.

Analizele au fost efectuate pentru acele combinații de valori numerice ale operanzilor de intrare care au reprezentat cazurile cele mai defavorabile din punct de vedere al propagării transportului prin structurile proiectate. De asemenea au fost efectuate analize pentru a confirma corectitudinea realizării funcției specifice a fiecărui bloc. A fost utilizat simulatorul MicroSim, în ale cărui librării de componente au fost modificați timpii de propagare prin diversele porți logice utilizate în conformitate cu datele furnizate de [67], corespunzătoare unei tehnologii CMOS 0,5 μm , adoptată ca și etalon pentru a permite o mai facilă comparație cu circuite similare studiate în referințele bibliografice. Astfel, timpii de propagare prin diversele porți logice sau de transmisie utilizate au fost considerați după cum urmează: 0,35 ns pentru porți SAU-exclusiv neconvenționale, mai rapide decât cele clasice, 0,2 ns pentru porți ȘI-NU și SAU-NU cu două intrări, 0,1 ns pentru inversoare CMOS și pentru tranzistoare de trecere.

Deoarece nu am dispus de un soft profesional pentru simulare după extragere din layout, am încercat să modelez funcționarea reală a structurilor numerice proiectate prin introducerea unor întâzieri suplimentare pe anumite căi de procesare a informației. Astfel au fost folosiți arbori de inversoare CMOS pentru multiplicarea capacității de comandă a unor linii de selecție ce comandau un număr mare de grile MOS; de asemenea valoarea capacităților parazite ce se acumulează datorită cascaderii porților de transmisie din circuite selectoare consecutive a fost limitată prin intercalarea unor circuite latch, realizate tot cu inversoare CMOS.

Datorită numărului mare de porți logice utilizate la construcția unui nivel pipeline (de ordinul miilor sau chiar mai mult), nu întotdeauna simulatorul a făcut

față situației, caz în care simularea s-a făcut doar pe blocuri componente, nu și pe ansamblu. În acest caz, verificarea prin simulare a blocurilor componente a condus la un calcul acoperitor al timpului total de propagare printr-un nivel pipeline deoarece pentru o procesare dată, blocurile componente ale ansamblului pot să nu se găsească simultan în cel mai defavorabil caz.

4.1 Rezultatele simulărilor blocurilor și subblocurilor componente ale subunității de calcul logaritmice

Subunitatea logaritmică, concepută în prima variantă după metoda descrisă de F. S. Lai în [50], dar folosind algoritmul de conversie de format modificat (paragraful 2.2.2), a fost organizată pe șase niveluri pipeline (figura 2.14). Nivelurile 2 și 5 ce conțin multiplicatoare de 12×12 biți sunt nivelurile critice întrucât prezintă cel mai mare timp de propagare și anume 6,3 ns, confirmat în diagrama din figura 4.1.

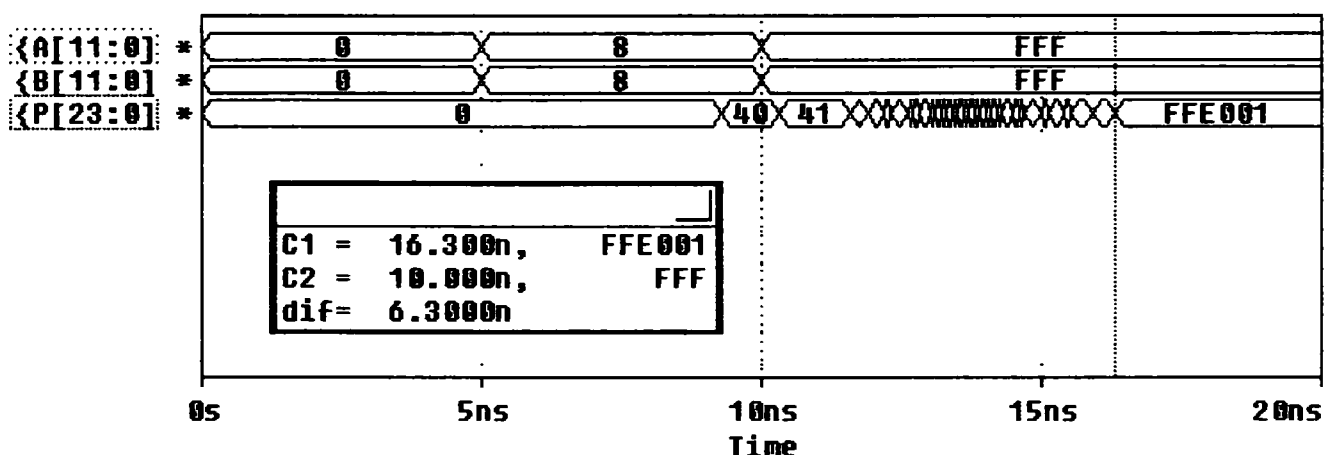


Fig. 4.1 Simularea multiplicatorului de 12×12 biți

Așa cum se observă în figura 4.1, dacă factorii A și B de 12 biți au valoarea maximă FFFh, atunci produsul P este reprezentat pe 24 de biți și are valoarea corectă FFE001h, obținută după 6,3 ns de la aplicarea stimulilor. Toate produsele parțiale având valoarea FFF, fiecare sumator complet din arborele Wallace al ariei de înmulțire se va găsi în cel mai defavorabil caz al său, în care toate cele 3 intrări ale sale au valoarea logică 1. Funcționarea corectă a multiplicatorului se poate constata mai ușor în cazul când operanzii au valoarea 8. Rezultatul este 64 adică 40h și e obținut într-un timp mai mic de 5 ns.

Următorul nivel pipeline consumator de timp al structurii din figura 2.14 este cel care conține blocul ALU. Pentru a fi obținuți timpi de propagare cât mai mici au fost modificate sau chiar concepute diferite sumatoare rapide pe 32, 16 și 8 biți. Astfel, dacă sumatorul CLA clasic cu structură modulară pe 32 de biți produce rezultatul UF pentru combinația cea mai defavorabilă de operanzi UD și UE (justificată în paragraful 2.2.2.4, pag.55, aliniatul 2) în 5,7 ns, așa cum se observă

în diagrama din figura 4.2.a, sumatorul CLA cu structura blocurilor componente modificată de autorul tezei produce același rezultat în doar 4,3 ns, așa cum se observă în diagrama din figura 4.2.b.

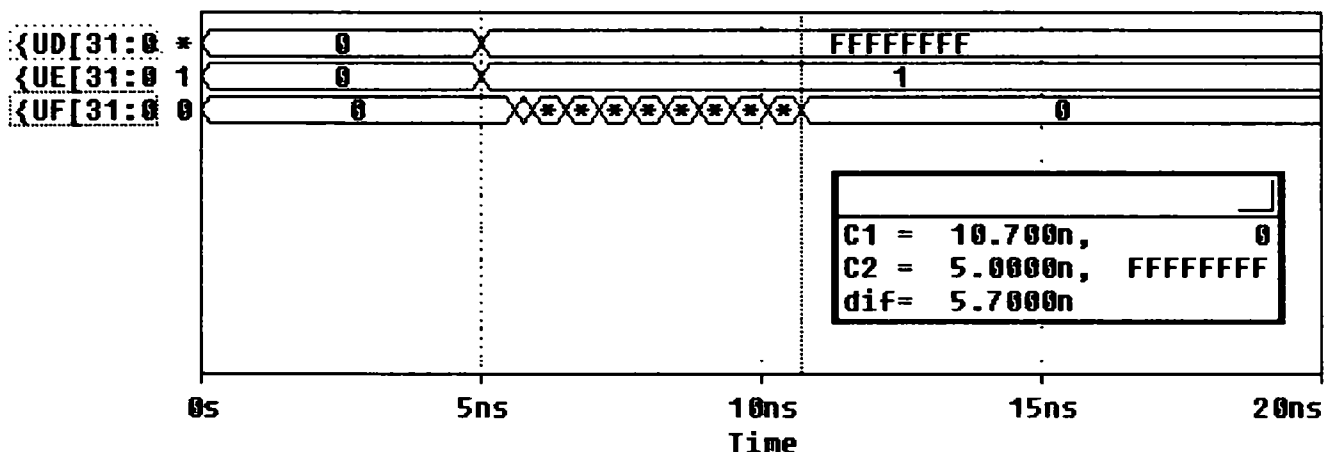


Fig. 4.2.a Simularea sumatorului CLA pe 32 de biți clasic

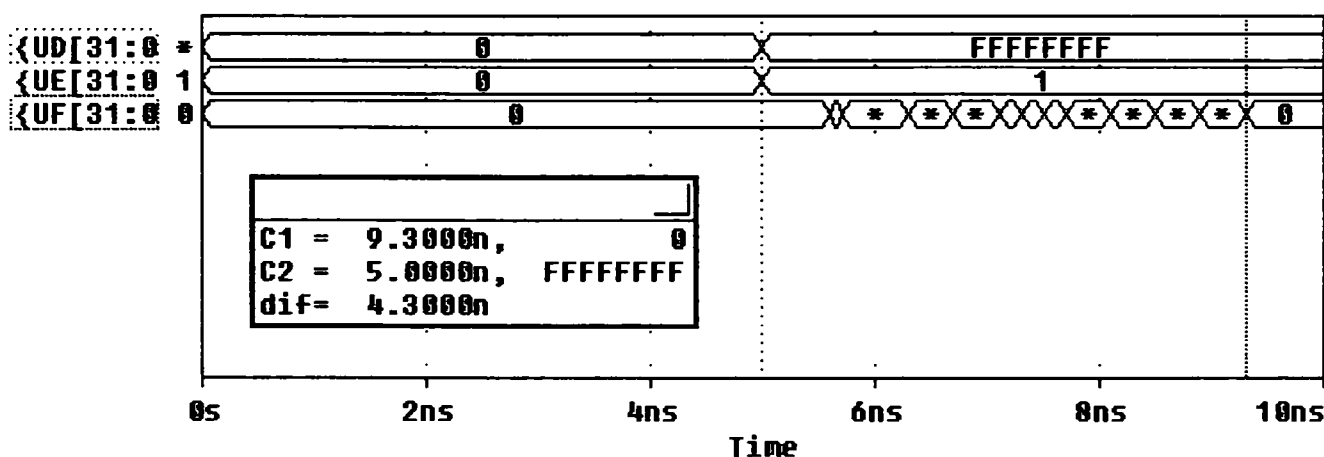


Fig. 4.2.b Simularea sumatorului CLA pe 32 de biți modificat

Pentru micșorarea în continuare a timpului de producere a rezultatului, a fost conceput un nou tip de sumator pe 32 de biți ce a avut ca punct de plecare observația importantă făcută în paragraful 2.2.2.4, pag. 55, aliniatul 3. Noul tip de sumator este în ansamblu un sumator cu selecția transportului, dar la care blocurile componente sunt sumatoare de tip CLA pe un număr mai mic de biți. Din acest motiv au fost analizate și sumatoare CLA pe 16 și 8 biți cu structura blocurilor componente modificată. Așa cum se observă în diagramele din figurile 4.3.a și 4.3.b, au fost obținuți timpi maximi de propagare de 3,6 ns și respectiv 2,5 ns. În realitate acești timpi sunt cu 0,1 ns mai mici deoarece faptul că transportul inițial al acestor sumatoare este 0 în cazul cel mai defavorabil (vezi paragraful 2.2.2.4) conduce la o simplificare a blocurilor de tip B' și B'' situate pe rama inferioară a structurii modulare a sumatorului CLA.

Structura optimă a noului sumator cu selecția transportului, din punct de vedere al vitezei de operare, a fost cea în care blocurile componente erau reprezentate de sumatoare CLA pe 8 biți. Această structură reprezintă un optim

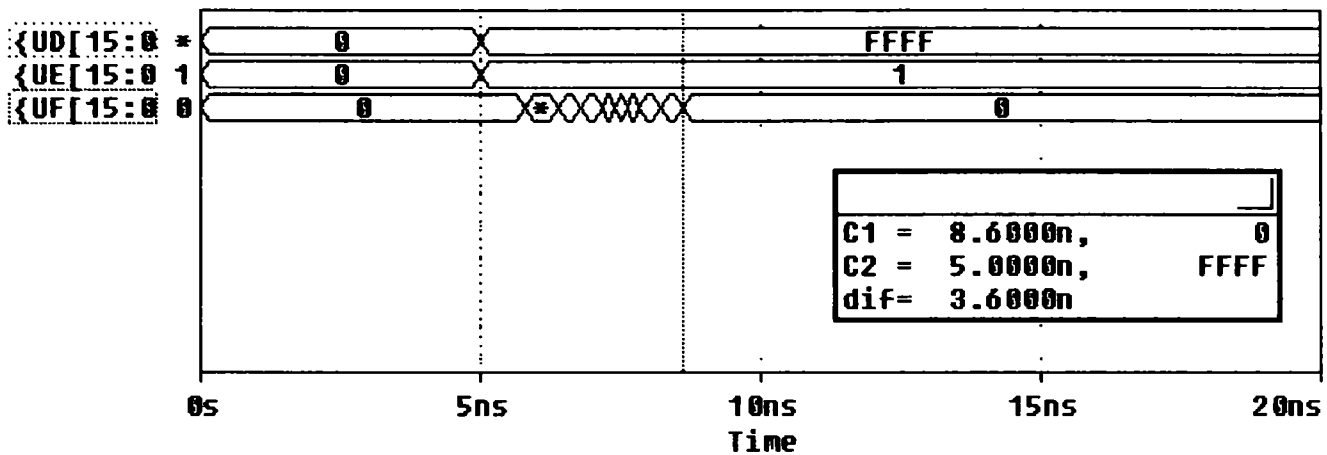


Fig. 4.3.a Simularea sumatorului CLA pe 16 biți modificat

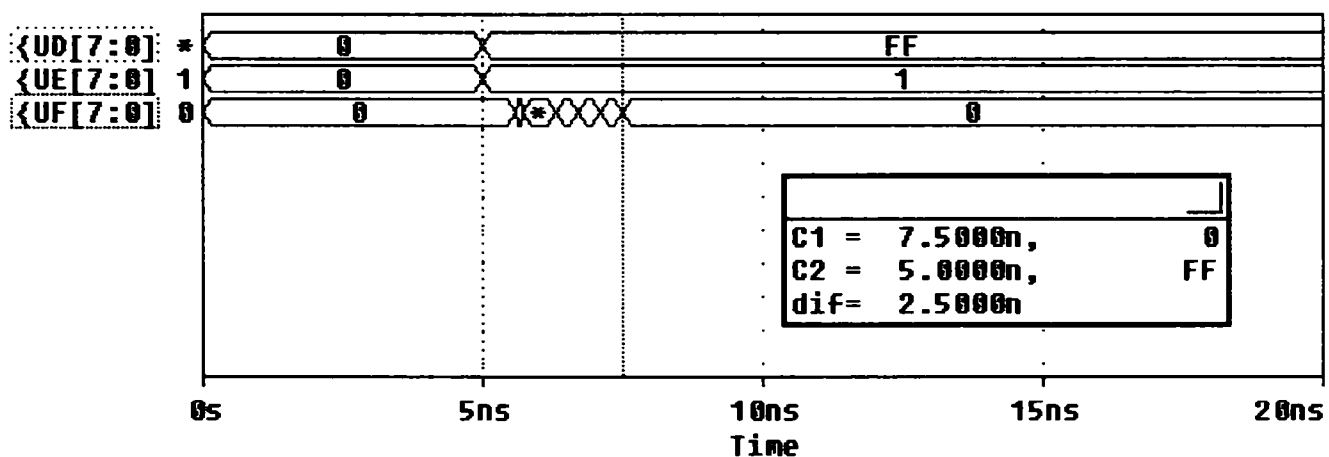


Fig. 4.3.b Simularea sumatorului CLA pe 8 biți modificat

deoarece intervalul de timp cuprins între generarea transportului furnizat de sumatorul CLA de rangul cel mai mic și stabilirea celor mai puțin semnificativi 8 biți ai sumei acoperă aproape în întregime timpul necesar pentru multiplicarea capabilității de comandă a liniilor de selecție a transportului de intrare corect pentru sumatoarele CLA de rang superior. Din acest motiv, timpul de producere al rezultatului complet este cu doar 0,2 ns mai mare decât al unui sumator CLA pe 8 biți. Așa cum se observă în diagrama din figura 4.4, acest timp este de 2,6 ns.

Combinăția de operanzi ce conduce la cel mai defavorabil caz din punct de vedere al propagării transportului prin circuit, rezultată în urma împărțirii cuvintelor ce se însumează în patru secțiuni de câte 8 biți, a fost găsită pe baza următoarelor constatări:

- *sumatorul CLA de rangul cel mai mare (al cărui rezultat se stabilește cel mai târziu) trebuie să se găsească în cel mai defavorabil caz al său ceea ce corespunde aplicării la intrarea sa a combinației FFh și 01h;*
- *sumatoarele CLA care funcționează cu transport inițial 0 pot reprezenta o cale critică de propagare, întrucât asigurarea transportului inițial 1 nu se regăsește în nici o combinație de intrare ce reprezintă cel mai defavorabil caz (vezi aliniatul 2, pag. 55) pentru sumatorul CLA. Pe de altă parte, aceste sumatoare nu trebuie să*

genereze transport 1, pentru a se păstra aceeași situație defavorabilă la sumatorul de rang imediat superior. În consecință cel mai defavorabil caz corespunde combinației de cuvinte de intrare 7Fh și 01h. În compensație, liniile de selecție a rezultatului corect de la fiecare pereche de sumatoare de același rang, rămân tot timpul pe 0 și selecția este foarte rapidă. Așadar, structura particulară a noului sumator propus de autorul tezei nu permite aflarea simultană în cazul cel mai defavorabil a sumatoarelor componente;

- sumatoarele CLA care funcționează cu transport inițial 1 pot reprezenta o cale critică de propagare, deși nu se încadrează în cel mai defavorabil caz al sumatorului de tip CLA, deoarece generarea transportului de ieșire al sumatorului și deci selecția rezultatului corect de la perechea de sumatoare din aval se face cu întârziere. Așadar și în acest caz are loc o compensare benefică a întârzierilor de pe liniile de selecție și de pe circuitele sumatoare.

În concluzie, simularea se va face pentru seturile de operanzi FF7F7F7Fh și 01010101h, respectiv FFFFFFFFh și 00000001h, precum și pentru o combinație între acestea. Așa cum se observă în diagrama din figura 4.4, ambele seturi conduc la obținerea unui timp de propagare de 2,6 ns. Orice altă combinație între aceste două seturi conduce la obținerea mai rapidă a rezultatului.

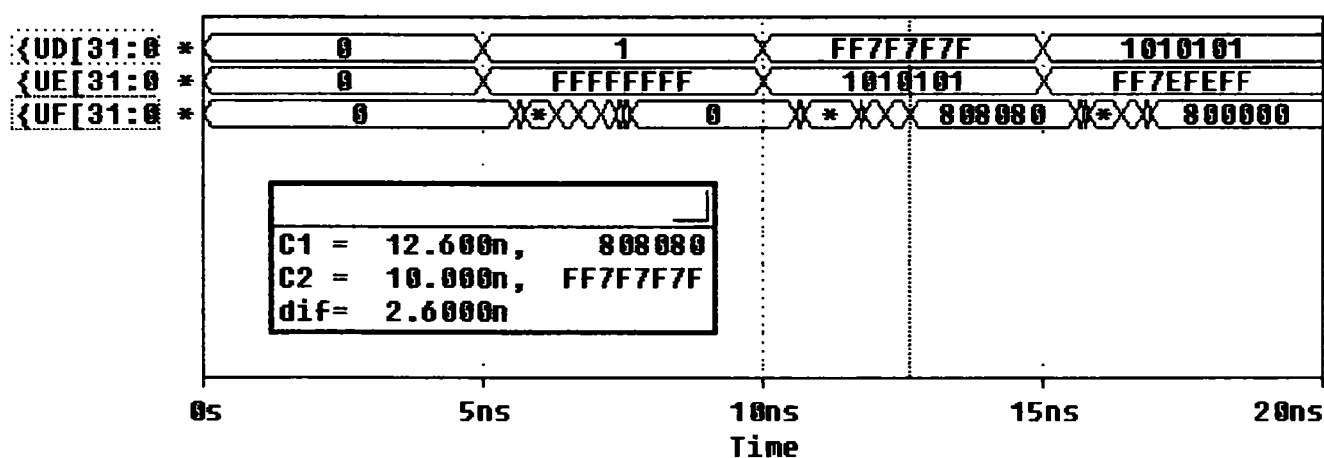


Fig. 4.4 Simularea noului sumator rapid pe 32 de biți

Etapă următoare a constat în simularea blocului ALU din figura 2.14, adică a circuitului sumator/scăzător prezentat în figura 2.9. Deoarece acest circuit a fost construit cu două sumatoare de tipul celui a cărui simulare a fost prezentată anterior, combinația cea mai defavorabilă de operanzi coincide cu cea din situația anterioară, dacă se efectuează o operație de adunare (/SEL=1, figura 2.9). Dacă se efectuează însă o operație de scădere, cazul cel mai defavorabil se obține pentru rezultat negativ (pentru a fi selectat rezultatul de la blocul Sumator1 cu transport inițial 0, care funcționează mai lent), atunci când descăzutul și scăzătorul iau cele mai mari valori posibile, diferite între ele. Așa cum se observă în diagrama din figura 4.5, dacă se efectuează scădere (/SEL=0), rezultatul UC are valoarea corectă “-1”, adică 80000001h cu bitul de semn inclus, după un timp de 3,3 ns. Și operația de adunare se efectuează tot în 3,3 ns, în cazul cel mai defavorabil, deoarece toate

blocurile selectoare din circuitul din figura 2.9 (Inversor1, Inversor2 și Selector) sunt incluse în calea critică de propagare.

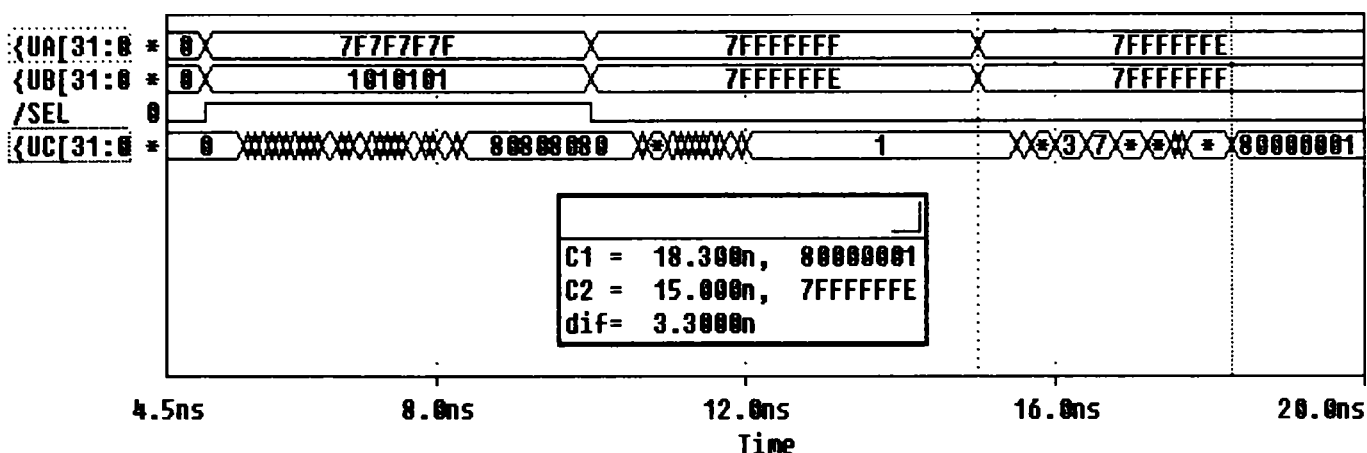


Fig. 4.5 Simularea circuitului de adunare/scădere din componența blocului ALU

Deoarece însă la blocul Inversor1 linia “/Sel” este deja setată (odată cu stabilirea tipului operației pe care o execută subunitatea logaritmică) înainte ca circuitele de logaritmare din figura 2.3 să furnizeze cei doi operanzi pentru ALU, întârzierea produsă la multiplicarea capacității de comandă a liniei /Sel, de 0,5 ns, este eliminată. În concluzie, timpul maxim de operare pentru ALU este de 2,8 ns.

Tot în nivelul 3 al structurii pipeline din figura 2.14 mai există blocurile “Sumator CSA+CLA”. Așa cum s-a precizat în paragraful 2.2.2.1, pag. 47, timpul maxim de propagare printr-un bloc CSA este dat de timpul de propagare prin două porți SAU-exclusiv cascade, adică 0,7ns. Timpul maxim de propagare prin sumatorul pe 23 de biți din componența aceluiași bloc, în cazul cel mai defavorabil este de 2,5 ns, așa cum se observă în diagrama din figura 4.6.

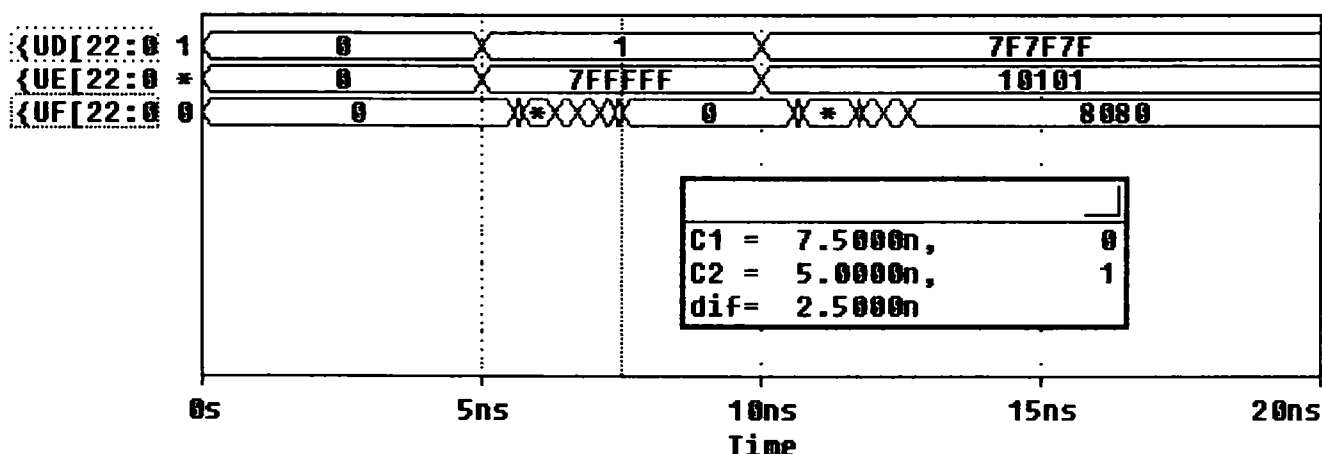


Fig. 4.6 Simularea sumatorului pe 23 de biți din blocul “Sumator CSA+CLA”

În concluzie, timpul total de propagare prin nivelul pipeline 3 al structurii din figura 2.14 este $0,7+2,5+2,8 = 6$ ns, așa cum s-a precizat în paragraful 2.2.2.5. Astfel, toate blocurile sau subblocurile componente ale structurii pipeline din figura 2.14 au fost simulate, confirmând rezultatele prezentate în paragraful 2.2.2.

În continuare vor fi prezentate rezultatele simulărilor blocurilor și subblocurilor componente ale subunității logaritmice proiectate prin implementarea metodei logaritmului mascat și a cărei schemă bloc a fost organizată tot pe 6 niveluri pipeline (figura 2.25). În acest caz au fost folosite aceleași tipuri de sumatoare performante de concepție proprie, însă extinse pe un număr mai mare de biți, după cum s-a precizat în paragraful 2.2.3.3. În plus, au fost folosite blocuri de compresori 4:2 pentru care timpul de propagare este dat practic de timpul de propagare printr-un singur compresor, așa cum s-a explicat în paragraful 2.2.3.2. Compresorul 4:2 este astfel considerat celula de bază pentru structurile numerice din nivelurile pipeline 2 și 5, dar el este prezent și în blocul ALU de pe nivelul 3. Pentru a demonstra funcționarea corectă a compresorului 4:2 L propus de autorul tezei, au fost efectuate simulări pentru toate combinațiile posibile ale intrărilor circuitului. Utilizând aceiași stimuli, au fost efectuate apoi simulări și pentru compresorul 4:2 optimizat, prezentat în [72]. Rezultatele obținute sunt prezentate în diagramele din figurile 4.7.a și 4.7.b.

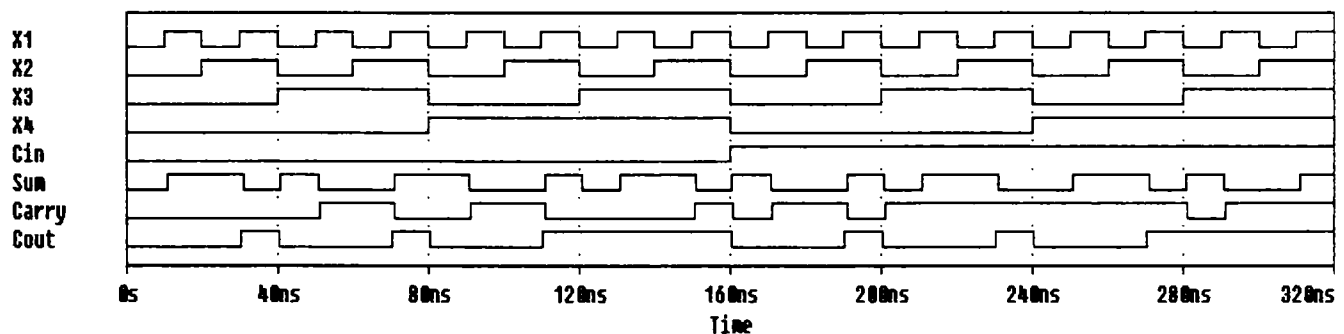


Fig. 4.7.a Simularea compresorului 4:2 L

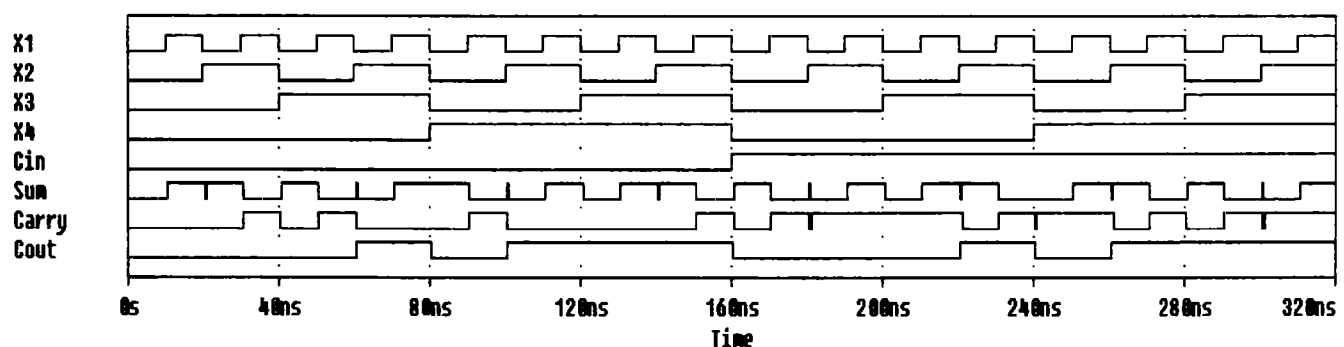


Fig.4.7.b Simularea compresorului 4:2 prezentat în [72]

După cum se observă în aceste diagrame, pe cele cinci linii de intrare X1, X2, X3, X4 și C_{in} au fost aplicate toate combinațiile posibile de valori logice. Analiza a fost efectuată în primă fază pe intervale de timp de câte 10 ns începând din momentul aplicării unei noi combinații la intrare, adică pe intervale de timp mult mai mari decât timpii de propagare prin circuit. Scopul acestei analize a fost de a verifica funcțiile logice asigurate de circuit, în conformitate cu definiția dată pentru compresorul 4:2 în paragraful 2.2.3.2. Se constată astfel că ieșirea S_{um} este

aceeași în ambele diagrame și în concordanță cu funcția a) enunțată la pag. 70. Ieșirile C_{arry} și C_{out} prezintă aceeași stare logică în ambele diagrame, dacă suma $X1+X2+X3+X4+C_{in}$ este egală cu 0, 1, sau 5, fapt care este conform cu funcțiile d) și f) ce trebuie asigurate de circuit (pag.70). Dacă suma menționată este egală cu 2 sau 3, atunci suma $C_{arry}+C_{out}$ este întotdeauna egală cu 1, în ambele diagrame; cu alte cuvinte ieșirile C_{arry} și C_{out} sunt complementare în ambele circuite. Acest fapt este în concordanță cu funcția e) a compresorului 4:2.

Pentru compresorul 4:2 L s-a făcut și o analiză în termeni de viteză pentru cea mai defavorabilă combinație de intrare care va determina propagarea unui transport pe calea critică de propagare a circuitului. Așa cum se observă în figura 2.23, calea critică de propagare include ieșirea C_{out} a compresorului din amonte (care poate genera un transport după 0,4 ns de la schimbarea stării logice a intrărilor, deci mai târziu decât o poartă SAU-exclusiv de pe intrarea compresorului curent), intrarea C_{in} a compresorului curent și ieșirea C_{arry} a acestuia. Una din combinațiile de intrare posibile ce permite generarea unui transport pe această cale este cea utilizată în simularea al cărei rezultat este prezentat în figura 4.8. Se observă că ieșirea C_{arry} devine stabilă după 1,1 ns de la aplicarea stimulilor. Această valoare reprezintă deci timpul maxim de propagare prin compresorul 4:2 L.

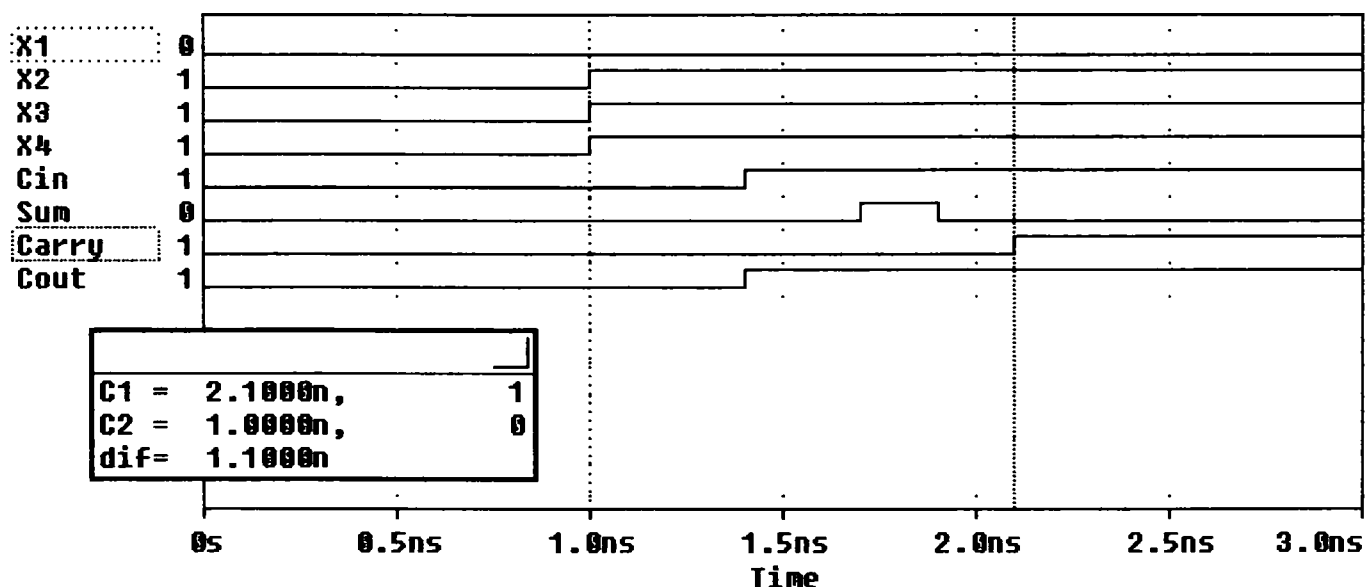


Fig. 4.8 Simularea compresorului 4:2 L în cazul cel mai defavorabil

Următoarea etapă a constat în simularea arborelui Wallace care ocupă întreg nivelul 2 al structurii pipeline în care a fost organizată subunitatea de calcul logaritmă, în cazul implementării metodei logaritmului mascat. Deoarece la construcția acestui nivel au fost utilizate peste 3000 de porți logice, porți de transmisie și inversoare CMOS, simulatorul de care am dispus nu a mai făcut față volumului imens de calcul, motiv pentru care am apelat la un artificiu. Acesta a constat în secționarea uneia dintre magistralele interne de date, ale cărei terminale au fost extrase la două porturi exterioare, cel de ieșire notat TT[13:0], iar cel de

intrare notat $T[13:0]$. Astfel, simularea arborelui Wallace s-a efectuat în două etape: în prima etapă a fost găsită valoarea numerică de pe liniile magistralei de ieșire $TT[13:0]$ și timpul scurs din momentul aplicării stimulilor până când această valoare a devenit stabilă; în a doua etapă au fost aplicați aceeași stimuli plus încă unul la portul de intrare $T[13:0]$ ce introducea valoarea numerică găsită anterior și întârzierea aferentă.

Simularea a fost făcută în scopul verificării funcționării corecte a circuitului. În ce privește timpul de propagare în cazul cel mai defavorabil, așa cum s-a arătat în paragraful 2.2.3.2, pag.69 și paragraful 2.2.3.3, pag. 74, acesta este de 4 ns. După efectuarea simulării, s-a constatat că întârzierea obținută în urma propagării informației binare prin structurile hard ce formează al doilea nivel pipeline s-a încadrat în 4 ns, așa cum se observă în diagramele din figura 4.9.

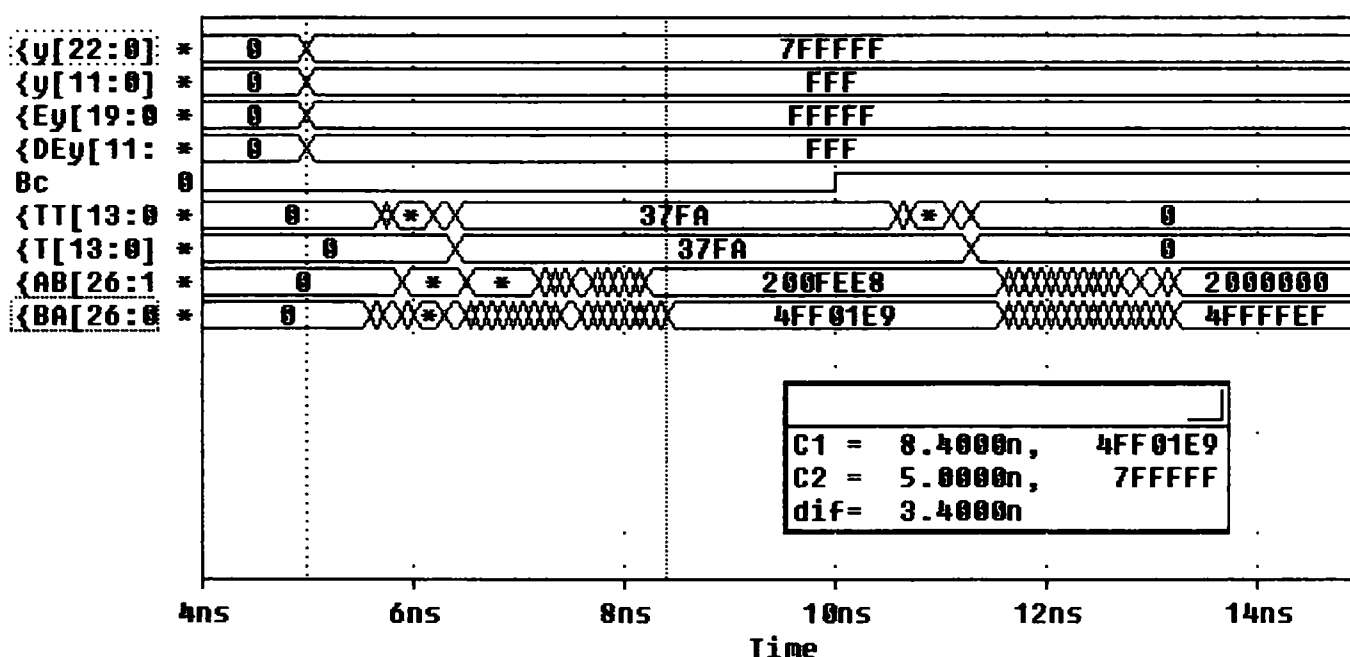


Fig. 4.9 Simularea structurii hard a nivelului pipeline 2 al subunității logaritmice concepută prin implementarea metodei logaritmului mascat

Așa cum se observă în figură, datele de intrare y , y_2 ($y[11:0]$), Ey și DEy au valorile maxime, pentru a încadra primul etaj al arborelui Wallace în cazul cel mai defavorabil, dar și pentru o verificare matematică mai ușoară a rezultatului final, ce s-a dovedit a fi corect. Cuvintele furnizate la ieșire, AB și BA pe 26 și respectiv 27 de biți sunt decalate cu un bit, datorită mecanismului prin care au fost produse, astfel că data AB va primi la dreapta un zero care va deveni LSB și va avea rangul -27 ca și LSB al lui BA . Cuvintele AB și BA sunt cele care trebuie “să mascheze” partea fracționară a logaritmului operandului aplicat la o intrare a subunității logaritmice, dacă valorile de intrare Ey și DEy ar fi furnizate de memoriile ROMA și ROMA'. Linia Bc reprezintă bitul de control memorat în locațiile lui ROMA și asigură funcționarea algoritmului descris în paragraful 2.2.3.1.

Structura hard a nivelului pipeline 5 este identică cu cea simulată anterior, diferențele regăsindu-se doar în valorile datelor aplicate la intrarea sa (v. pag.67).

În continuare va fi prezentată simularea unei secțiuni importante a blocului ALU de pe nivelul pipeline 3 al subunității logaritmice. Această secțiune cuprinde circa 1000 de porți logice ce reprezintă aproximativ 60% din numărul total de porți utilizate la construcția blocului ALU. Simularea întregului bloc nu a fost posibilă datorită limitărilor impuse de simulator, dar acest lucru nu a împiedicat asupra evaluării și verificării timpului total de propagare prin nivelul pipeline 3 întrucât blocurile care nu au fost incluse în secțiunea simulată prezintă un timp de propagare cunoscut și verificat prin simulări anterioare. Este vorba de blocul de compresori (1,1 ns - timp maxim de propagare) din figura 2.24 și cele două blocuri multiplexoare InversorB1 și InversorB2 care operează în paralel (0.2 ns - timp maxim de propagare). În consecință secțiunea simulată cuprinde blocurile principale Sumator1 și Sumator2 pe 36 de biți, precum și blocurile multiplexoare Inversor2 și Selector împreună cu poarta SAU-NU (fig. 2.24). Pentru a putea fi efectuată simularea, datele aplicate la intrare au fost distribuite pe câte două magistrale întrucât simulatorul nu a suportat stimuli extinși pe mai mult de 32 de biți. Deoarece sumatoarele utilizate sunt de același tip cu cel simulat în figura 4.5 a fost folosită aceeași structură a stimulilor pentru cazul cel mai defavorabil. Datorită faptului că un cuvânt de date din cele două furnizate de blocul de compresori nu conține bit în poziția LSB, unul dintre stimulii aplicați circuitului este reprezentat doar pe 31 de biți. Acest lucru constituie un avantaj deoarece combinația cea mai dezavantajoasă de operanzi la intrările sumatoarelor nu mai este posibilă. Au fost efectuate atât simulări corespunzătoare efectuării operației de adunare (Sop=0), cât și operației de scădere. Așa cum se observă în figura 4.10, timpul maxim de propagare a fost găsit de 2,4 ns. Cumulând la această valoare încă 0,2 ns și 1,1ns se obțin 3,7 ns ceea ce asigură o rezervă de 0,3 ns în raport cu perioada de tact impusă sistemului.

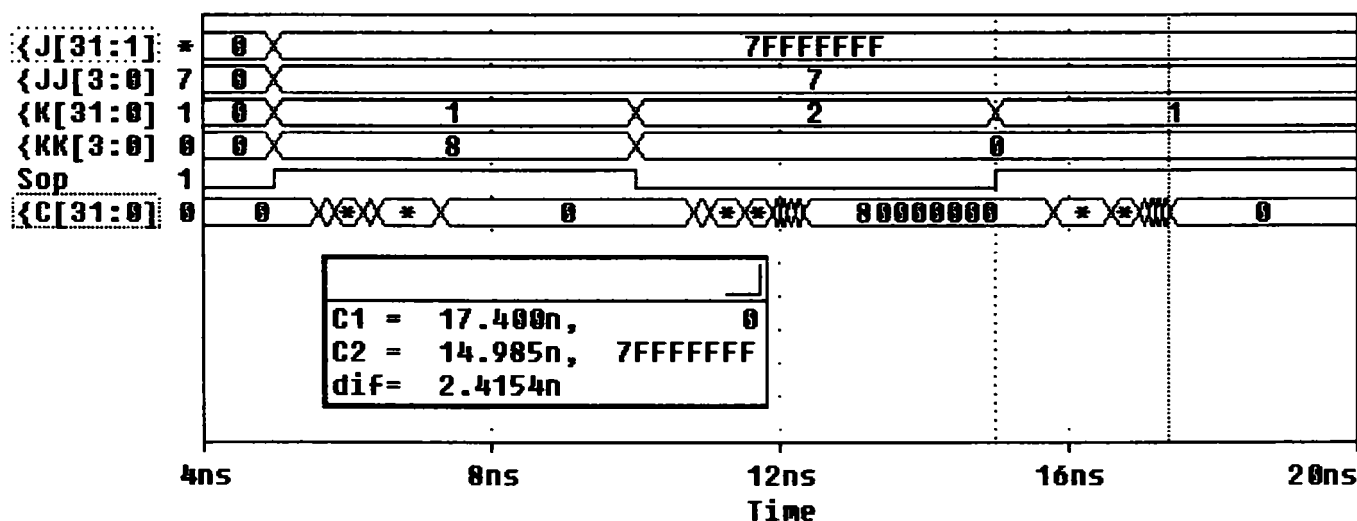


Fig. 4.10 Simularea secțiunii principale a blocului ALU

Ultimul nivel pipeline a cărui funcționare mai trebuie confirmată prin simulare este nivelul 6, în care se vor însuma două cuvinte de 27 de biți, din al

căror rezultat vor fi reținuți doar cei mai semnificativi 23 de biți pentru a forma mantisa din formatul în virgulă flotantă simplă precizie. Și în acest caz combinația cea mai dezavantajoasă de operanzi la intrările sumatorului nu mai este posibilă, deoarece operanzii sunt furnizați de un bloc de compresori, ca și în cazul discutat anterior. Din acest motiv, dar și pentru că se operează pe 27 de biți în loc de 32, a fost posibilă utilizarea cu succes a unui sumator CLA, de tipul celui din figura 2.11, mai puțin expansiv decât sumatoarele din ALU. În urma simulării a fost obținut un timp maxim de propagare de 3,6 ns, așa cum se poate constata în diagrama din figura 4.11.

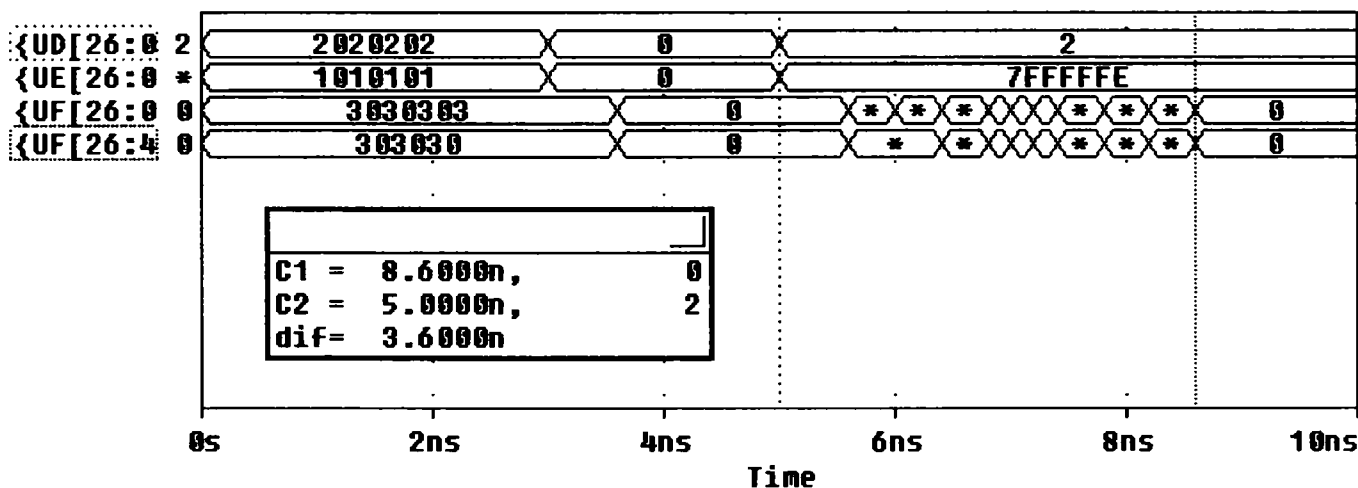


Fig. 4.11 Simularea sumatorului CLA pe 27 de biți de pe nivelul pipeline 5

4.2 Rezultatele simulărilor blocurilor și subblocurilor componente ale subunității de calcul în virgulă flotantă

La simularea celor trei variante de circuite de aliniere a mantiselor s-a avut în vedere atât determinarea timpului maxim de propagare prin circuit utilizând combinația cea mai defavorabilă de operanzi la intrare, cât și verificarea funcționării corecte a circuitului, pentru diverse combinații de operanzi.

Cuvintele de date ce se aplică acestui circuit sunt exponenții și mantisele celor doi operanzi, însă cazul cel mai defavorabil depinde doar de combinația de exponenți, întrucât mantisele s-ar propaga deplasate sau nedepasate în același timp prin selectoarele circuitelor “barrel shifter”, dacă liniile de selecție ale acestora ar avea valori logice inițial cunoscute și stabile. Deoarece în cazul **variantei 1** circuitele care efectuează cele două diferențe ale exponenților sunt de același tip cu sumatoarele din blocul ALU al subunității logaritmice, combinațiile cele mai defavorabile de exponenți vor avea aceeași structură. S-a ținut cont și de faptul că ambele sumatoare funcționează cu transport inițial 1. Combinația vizată corespunde setului de valori F7h și EFh, respectiv FFh și FFh pentru care au fost afișate întârzierile la propagarea celor două mantise în diagramele din figura 4.12. Așa cum se observă în diagrame, timpul maxim de propagare este de 2,4 ns în ambele situații.

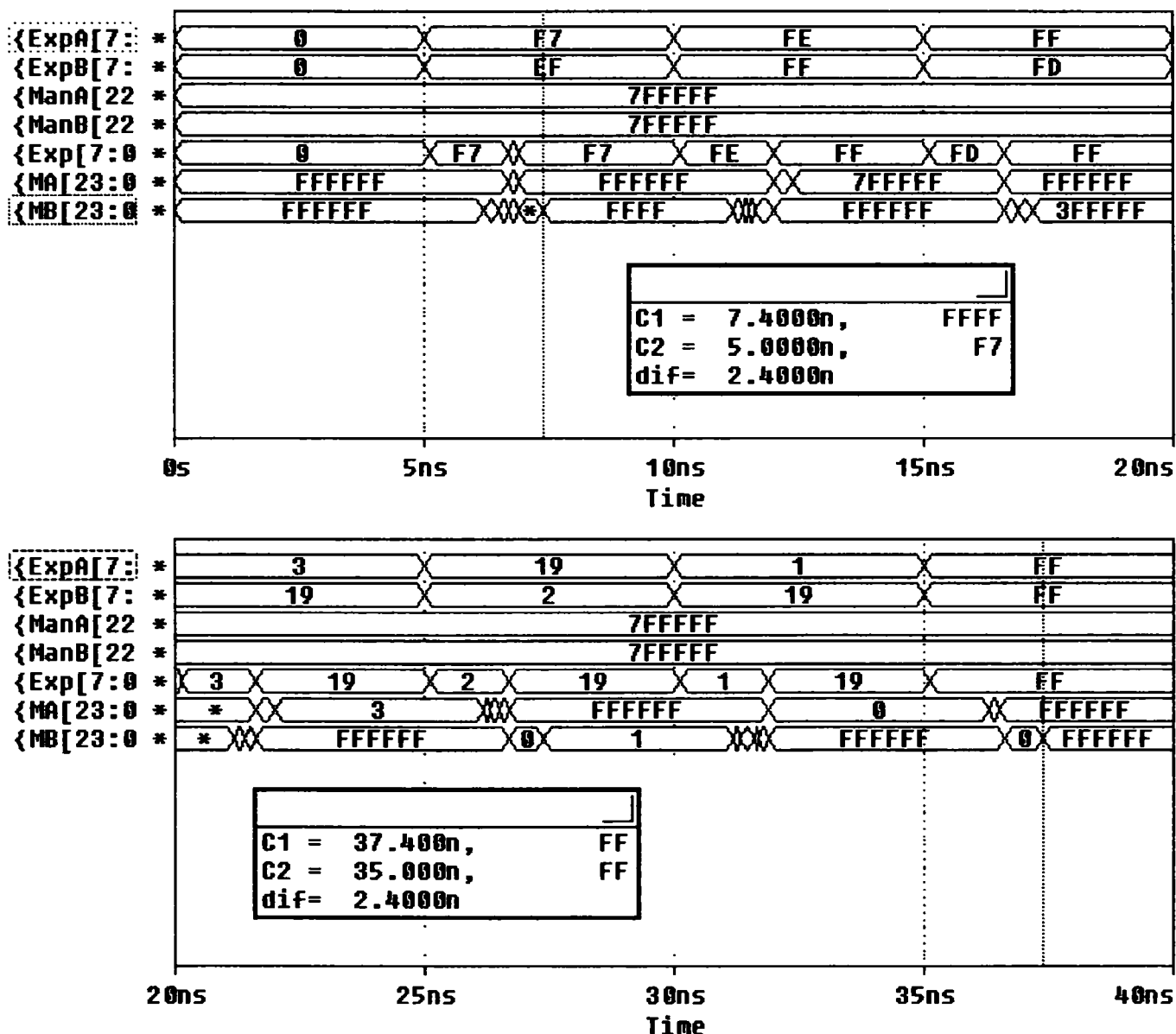


Fig. 4.12 Simularea circuitului de deplasare a mantiselor, varianta 1

În scopul verificării ușoare a deplasării corecte a mantiselor, acestea au fost stabilite ca având toți biții egali cu 1. De asemenea, pentru a pune în evidență deplasarea mantisei corespunzătoare exponentului mai mic situația relativă a celor doi exponenți a fost alternativ schimbată, la fiecare 5 ns. Așa cum se observă în figură, au fost obținute deplasări la dreapta ale mantisei inferioare cu 8, 1 și 2 biți în diagrama de sus și cu 22, 23 și 24 de biți în diagrama de jos. Deplasării cu 24 de biți la dreapta a mantisei inferioare îi corespunde evident anularea biților acesteia.

Se mai observă că la ieșire este întotdeauna selectat exponentul mai mare împreună cu mantisa asociată, nedeplasată.

La simularea **variantei 2** a circuitului de aliniere a mantiselor a fost folosit același fișier de stimuli din motivele practice care au fost explicate anterior. Totuși în acest caz combinația cea mai defavorabilă de exponenți nu mai este F7h și EFh deoarece s-a schimbat și tipul sumatoarelor componente ale sumatorului cu selecția transportului folosit pentru calcularea diferenței exponenților. Fiind folosite

sumatoare complete pe 1 bit cascade, combinația cea mai defavorabilă de exponenți este FFh și FEh ce corespunde calculării sumei FFh+01h+1, ce reprezintă cel mai defavorabil caz de la sumatorul simplu. Această combinație este de asemenea furnizată de fișierul de stimuli folosit. Ea conduce la obținerea unei întârzieri maxime de 3,7 ns așa cum se observă în diagramele din figura 4.13.

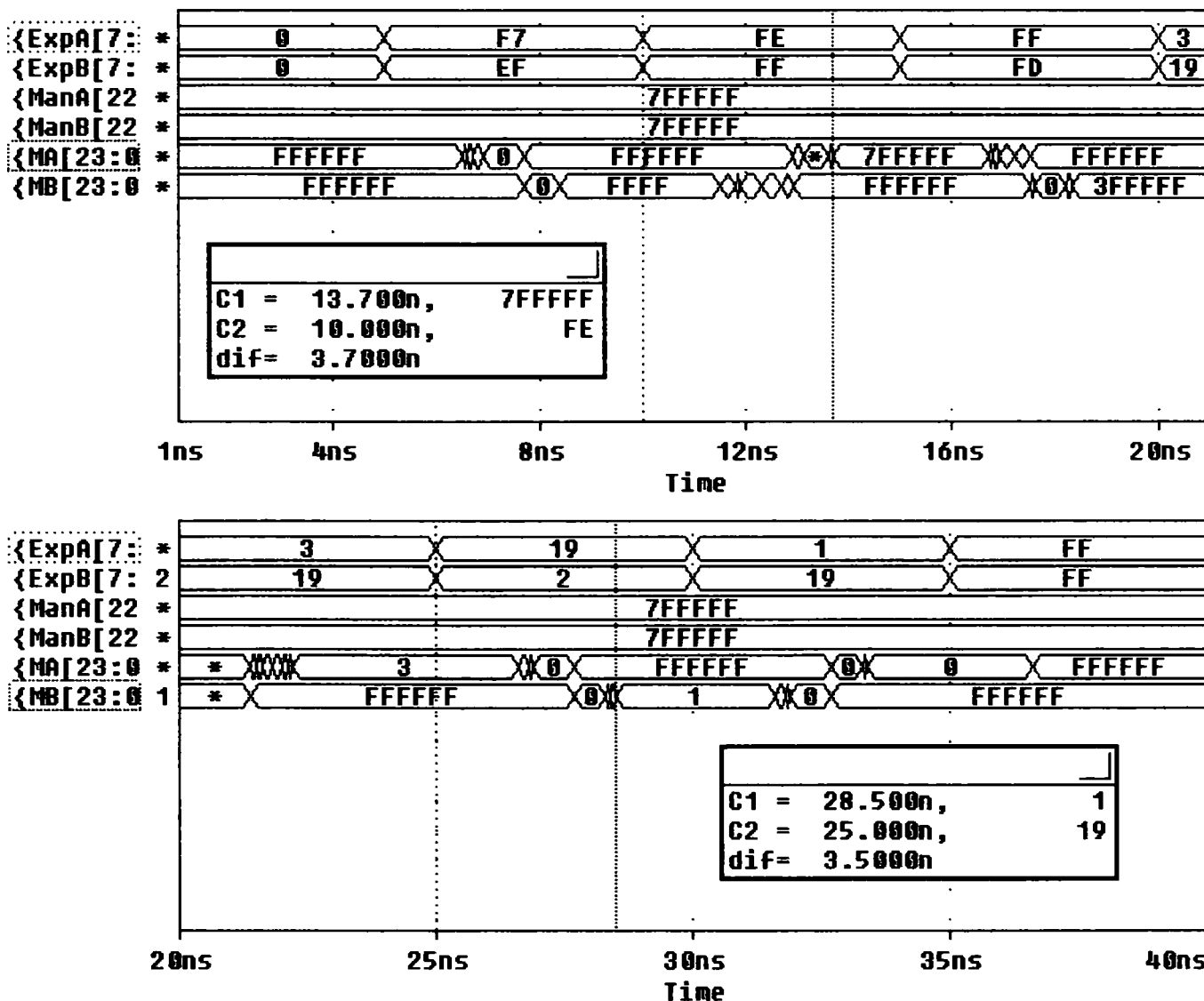


Fig. 4.13 Simularea circuitului de deplasare a mantiselor, varianta 2

De asemenea se constată că și acest circuit funcționează corect, iar întârzierea maximă produsă este sub 4 ns. El este deci mai potrivit a fi utilizat decât cel anterior deoarece ocupă o suprafață mai mică pe cipul de siliciu și întrunește condițiile de viteză impuse.

Elaborarea **variantei 3** a circuitului de deplasare a mantiselor, a condus însă la obținerea optimului, în sensul că a permis cea mai semnificativă economie de suprafață comparativ cu prima variantă (aproape 50%), iar întârzierea maximă obținută a fost de 4 ns, așa cum se observă în diagrama din figura 4.14.

Deoarece s-a utilizat același tip de circuit pentru efectuarea celor două diferențe ale exponenților, ca și în cazul variantei 2, combinația cea mai defavorabilă este identică cu cea din acest caz.

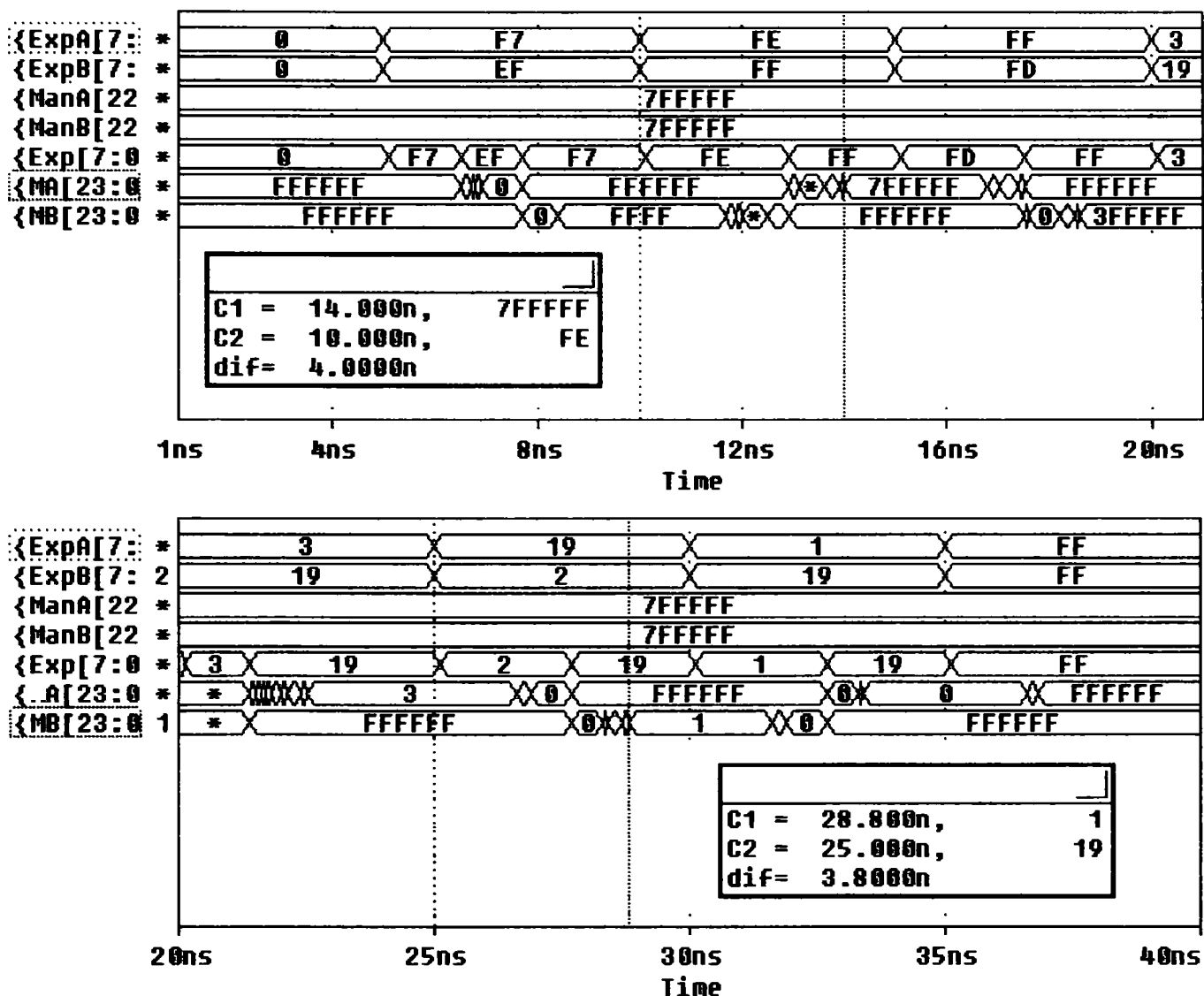


Fig. 4.14 Simularea circuitului de deplasare a mantiselor, varianta 3

Din diagramele prezentate în figura 4.14 se poate constata funcționarea corectă a circuitului. A fost folosit același fișier de stimuli ca și în simulările anterioare.

În continuare au fost efectuate simulări pentru stabilirea variantei optime de circuit sumator/scăzător cu care este realizat blocul ALU din subunitatea de calcul în virgulă flotantă, a cărei schemă bloc a fost prezentată în figura 3.6. Viteza acestui circuit depinde în principal de tipul celor două sumatoare pe 25 de biți din componența sa. În **varianta 1** de circuit sumator/scăzător, se utilizează sumatoare performante a căror schemă bloc a fost descrisă în figura 3.8 și sunt similare cu cele utilizate în blocul ALU al subunității logaritmice. Din acest motiv structura stimulilor urmează același tipar. Datorită complexității sporite a acestui bloc verificările vor fi efectuate pentru toate cele 8 combinații posibile din punct de vedere al semnelor operanzilor și al operației care se execută. În cazul în care se efectuează adunări propriu-zise combinația cea mai defavorabilă este FF7F7Fh și 010101h, iar în cazul în care se efectuează scăderi propriu-zise, această combinație

este FF7F7Fh și FEFEFEh (/FEFEFEh=010101h, în care operatorul “/” este operatorul negație). Timpii care sunt obținuți prin simulare includ un interval de timp de 1,1 ns necesar pentru stabilirea valorii logice a unor linii de selecție de la intrarea circuitului sumator/scăzător. Valoarea acestor linii nu depinde decât de semnele celor doi operanzi precum și de semnul operației care se efectuează (0 pentru adunare și 1 pentru scădere). În consecință stabilirea valorii logice a acestor linii, cât și multiplicarea capacității lor de comandă poate fi realizată în nivelul pipeline anterior, dacă se utilizează tranzistoare dimensionate corespunzător în latch-urile corespondente dintre primul și al doilea nivel pipeline.

În urma simulării circuitului în cazurile cele mai defavorabile, a fost obținut un timp maxim de propagare de 4 ns, așa cum se poate constata în diagramele din figura 4.15. Având în vedere observația anterioară, timpul real ce poate fi obținut în cazul cel mai defavorabil este de 2,9 ns.

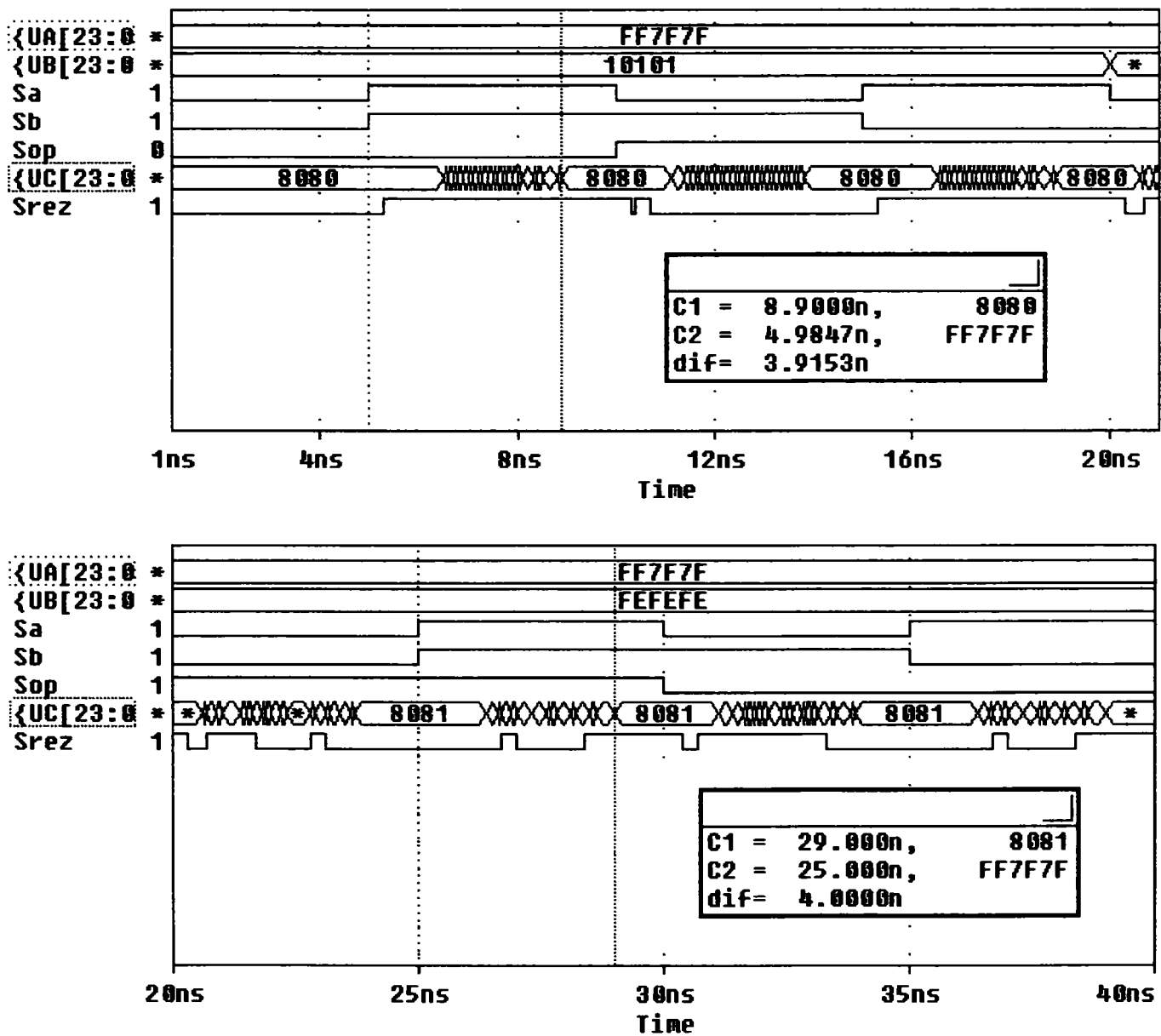


Fig. 4.15 Simularea variantei 1 a circuitului sumator/scăzător al subunității de calcul în virgulă flotantă

Stimulii aplicați circuitului au fost astfel construiți încât în primele patru intervale de câte 5 ns se efectuează adunări propriu-zise, iar în ultimele patru intervale se efectuează scăderi propriu-zise, în conformitate cu regulile enunțate în paragraful 3.2, pag. 86. Valoarea absolută a operanzilor este menținută aceeași, în schimb semnul lor se modifică periodic, ca și semnul operației, pentru a obține toate cazurile posibile de operare. Analizând diagramele din figura 4.15, se poate constata corectitudinea rezultatului obținut și a semnului acestuia, Srez. Se mai poate constata că bitul Srez este întotdeauna stabilizat înaintea rezultatului corespondent și că rămâne stabil până la următorul moment în care are loc schimbarea stimulilor aplicați la intrarea circuitului.

Cazurile cele mai defavorabile de la adunarea propriu-zisă, prezentate în diagrame, corespund însă unei situații de supradepășire mantisă, fapt indicat de bitul MSB de la blocul Sumator1 al circuitului din figura 3.6.

În **varianta 2** de circuit sumator/scăzător, am utilizat sumatoare CLA modificate, de tipul celui descris în figura 2.11, dar pe 25 de biți (inclusiv bitul de semn). Din acest motiv structura stimulilor se modifică corespunzător cu cel mai defavorabil caz de la acest tip de sumator. Verificările vor fi efectuate pentru toate cele 8 combinații posibile din punct de vedere al semnului operanzilor și al operației care se execută. În cazul în care se efectuează adunări propriu-zise combinația cea mai defavorabilă este FFFFFFFh și 000001h, iar în cazul în care se efectuează scăderi propriu-zise, această combinație este FFFFFFFh și FFFFFFFEh (/FFFFFFEh=000001h, în care operatorul “/” este operatorul negație). Din același motiv menționat anterior, timpul de propagare real prin blocul ALU din nivelul pipeline 2 este mai scurt cu 1,1 ns. Așa cum se observă în diagramele din figura 4.16, timpul maxim obținut prin simulare a fost de 5,1 ns, ceea ce corespunde unui timp real de propagare prin nivelul pipeline 2 al subunității în virgulă flotantă de 4 ns.

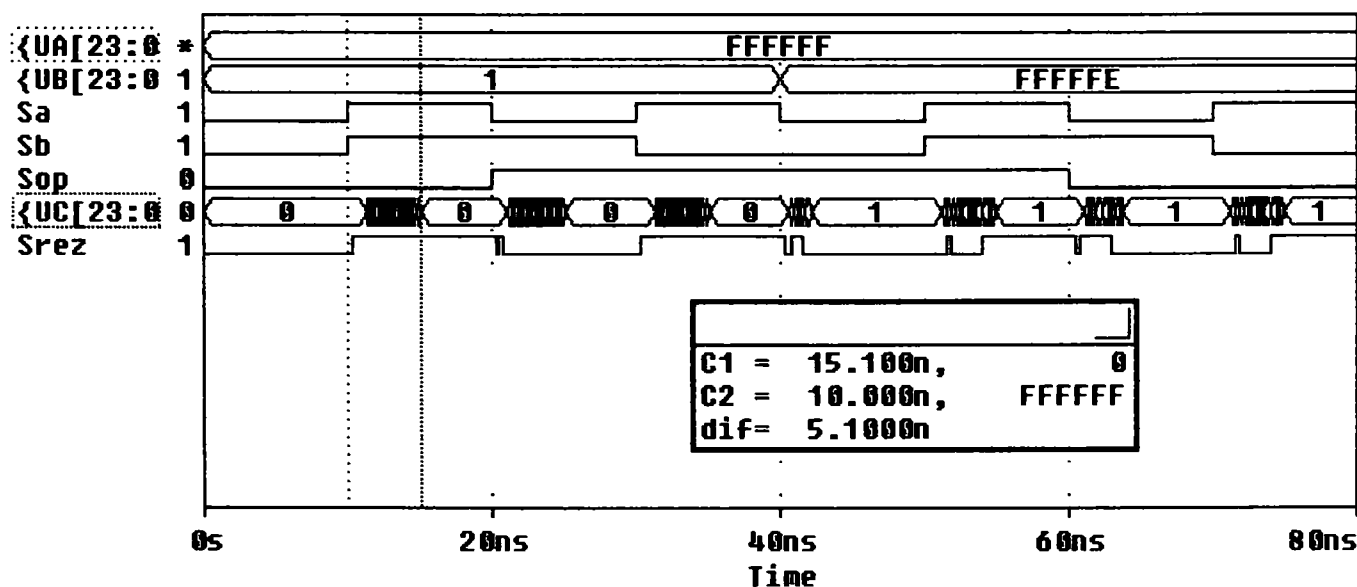


Fig. 4.16 Simularea variantei 2 a circuitului sumator/scăzător al subunității de calcul în virgulă flotantă

Analizând diagramele din figura 4.16, se poate constata corectitudinea rezultatului obținut și a semnelui acestuia, Srez. Cazul cel mai defavorabil corespunde unei adunări propriu-zise, în care este selectat întotdeauna rezultatul furnizat de blocul Sumator1 din figura 3.6 care funcționează cu transport inițial 0. Bitul MSB al acestuia va indica însă supradepășire mantisă, dar acest fapt nu influențează generarea corectă a bitului de semn al rezultatului.

În urma rezultatelor obținute la simulare, s-a constatat că varianta a doua de circuit de adunare/scădere reprezintă optimul întrucât este mai puțin expansivă în suprafață pe cip, iar timpul maxim de propagare pe care îl prezintă se încadrează în cele 4 ns impuse de subunitatea logaritmică.

În continuare vor fi prezentate rezultatele simulărilor celui de-al treilea nivel pipeline al subunității de calcul în virgulă flotantă. În figura 3.17 este prezentat rezultatul simulării blocului NZS din figura 3.9 care efectuează contorizarea zerourilor semnificative ale rezultatului furnizat de ALU în cazul efectuării unei operații de scădere.

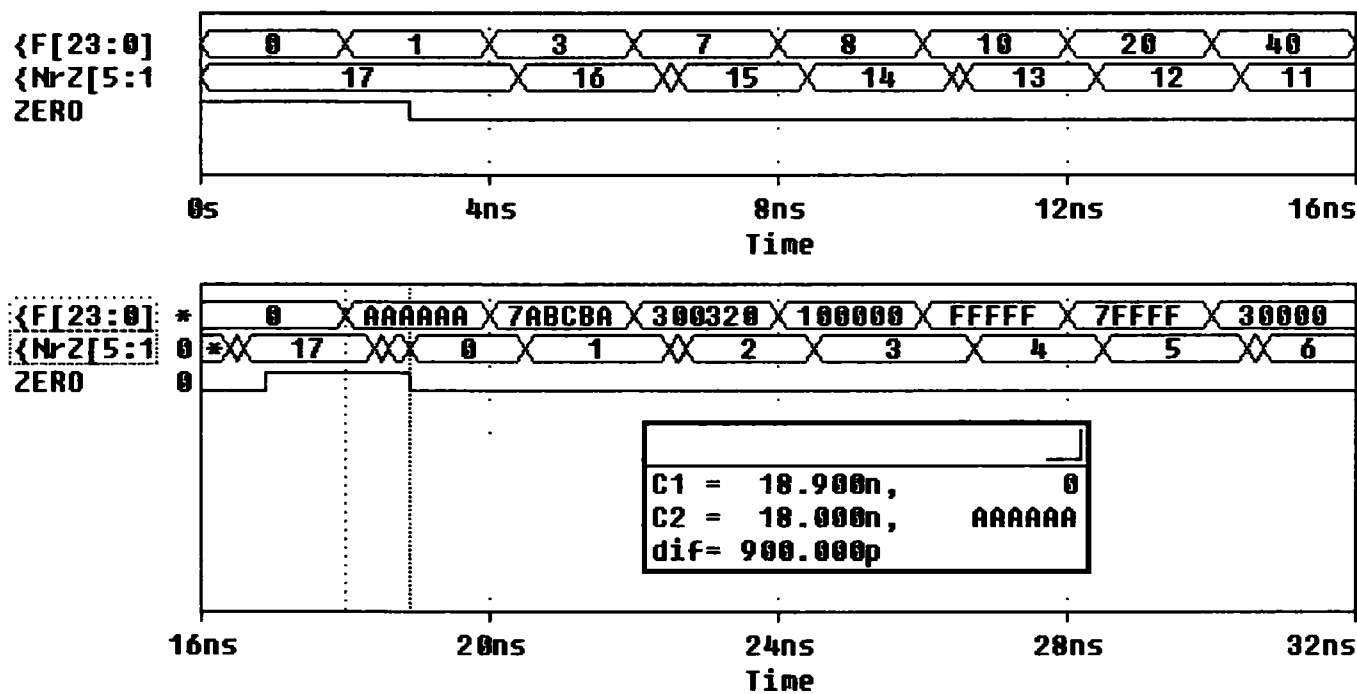


Fig. 4.17 Simularea numărătorului zerourilor semnificative ale diferenței mantiselor

După cum se observă în aceste diagrame, au fost efectuate simulări pentru cuvinte de date ce conțin 23÷17 (17h÷11h) și respectiv 0÷6 (00h÷06h) zerouri semnificative, rezultatul fiind redat corect în toate situațiile. În cazul în care toți biții diferenței mantiselor sunt egali cu zero, normalizarea nu mai are sens, acest fapt fiind indicat prin fanionul “zero” care va invalida rezultatul furnizat la ieșirea blocului NZS, care în acest caz ar genera valoarea maxim posibilă, adică 23 (17h).

Timpul maxim de propagare prin circuit a fost găsit de 0,9 ns. Diferența de 0,1ns față de valoarea evaluată teoretic în paragraful 3.3, pag.92, e dată de inversorul CMOS care a fost introdus pentru setarea pe 1 logic a fanionului “zero” atunci când se obține rezultat nul la ieșirea blocului ALU.

În continuare a fost simulat circuitul compus din blocurile BS (“barrel shifter”) și NRZ din figura 3.9 care deplasează rezultatul diferență furnizat de ALU, în scopul obținerii unei mantise normalizate. Deoarece ansamblul format din aceste două blocuri nu se găsește pe calea critică de propagare, așa cum s-a arătat în paragraful 3.3, simularea acestui circuit nu a urmărit decât verificarea funcționării sale corecte. În diagramele din figura 4.18 sunt prezentate rezultatele acestor simulări.

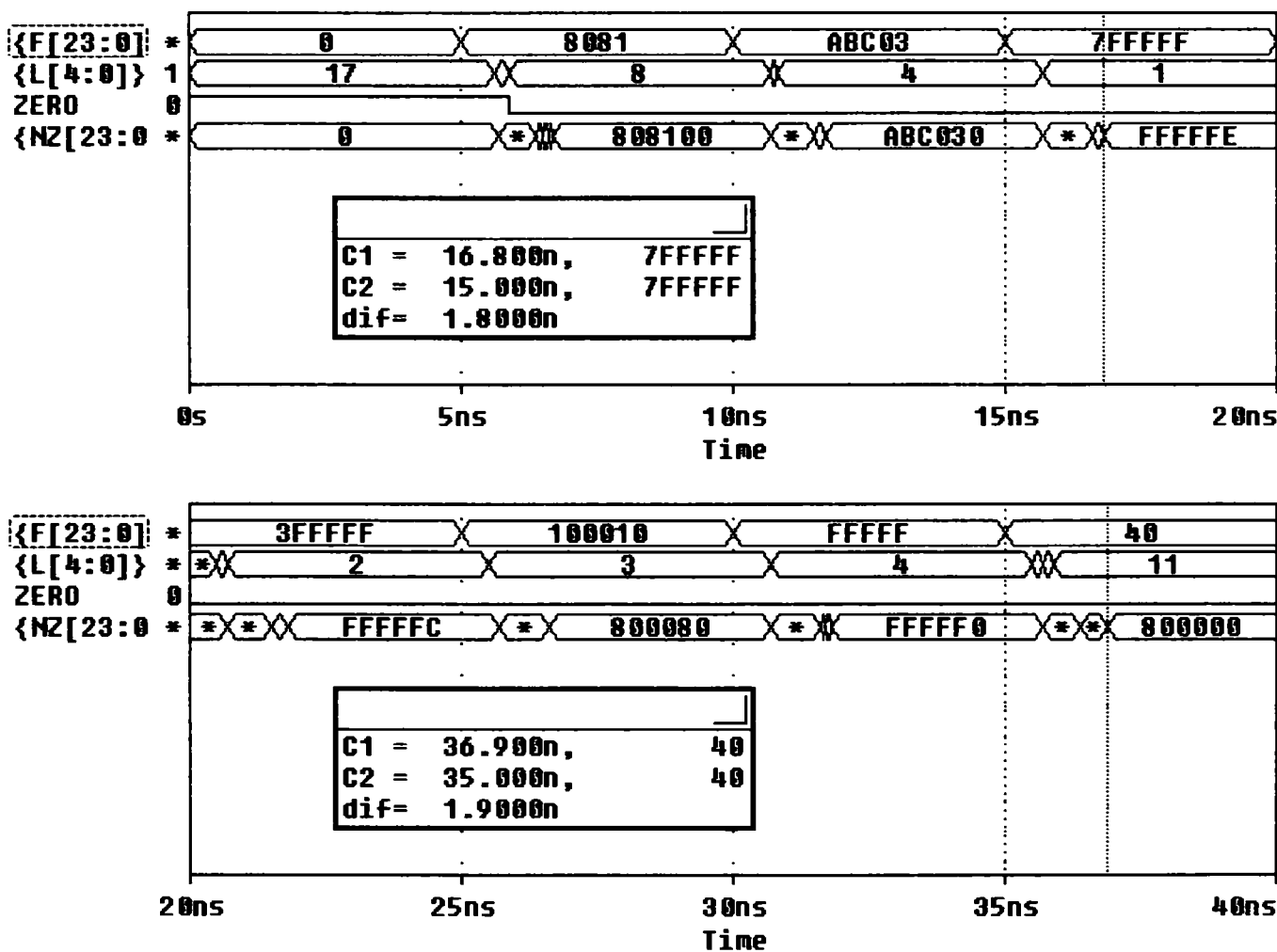


Fig. 4.18 Simularea circuitului de normalizare a rezultatului diferență furnizat de blocul ALU al subunității de calcul în virgulă flotantă

În această figură, F este rezultatul diferență furnizat de ALU, L este cuvântul pe 5 biți generat de blocul NZS, iar NZ reprezintă mantisa normalizată generată la ieșire. După cum se poate constata, toate datele numerice generate la ieșire au cel mai semnificativ bit egal cu 1 și poartă amprenta cuvintelor din care provin. În fiecare cuvânt de ieșire există un număr de zerouri în pozițiile cel mai puțin semnificative egal cu valoarea numerică a cuvântului generat de blocul NZS.

Întârzierea maximă obținută pentru cazurile aleatoare simulate este de 1,9 ns și se apropie de valoarea maximă estimată teoretic, de 2,1 ns, care oricum nu afectează întârzierea de pe calea critică de propagare din nivelul pipeline 3. Timpul maxim de propagare prin acest nivel se obține prin cumulara timpilor maximi obținuți în diagramele din figurile 4.17 (pentru blocul NZS) și 4.3.b (pentru blocul AJUST_EXP din figura 3.9), precum și a timpului de propagare prin multiplexorul INV/ZERO (fig. 3.9) de 0.2 ns. Se obține astfel un timp total de 3,6 ns, în condițiile în care suprafața necesită pentru implementarea pe cip este și în acest caz minimă.

Concluzii, contribuții și perspective

În lucrarea de față a fost prezentată o soluție proprie privind concepția unității de execuție numerică a unui procesor aritmetic performant, care implementează cele patru operații aritmetice elementare și care operează în simplă precizie. Pe baza analizelor efectuate în capitolul 1 al tezei cu privire la utilizarea metodelor numerice cunoscute de procesare a informației, s-a arătat la începutul capitolului 2 și ulterior s-a demonstrat practic prin proiectare propriu-zisă (capitolele 2 și 3) și s-a verificat prin simulare (capitolul 4) că se poate obține o viteză de operare foarte mare pentru un procesor aritmetic, atunci când se adoptă o structură hibridă pentru acesta. Noua structură a fost obținută prin includerea unei subunități de calcul logaritmice pentru efectuarea operațiilor de înmulțire și împărțire care prezintă intrări și ieșiri în format virgulă flotantă și a unei subunități de calcul în virgulă flotantă pentru efectuarea operațiilor de adunare și scădere.

Performanțele în ce privește viteza de operare, obținute prin utilizarea aritmeticii logaritmice, comparativ cu sistemele de calcul în virgulă flotantă sunt următoarele:

- deși o operație singulară de înmulțire se efectuează de 2,5÷3 ori mai lent decât una în virgulă flotantă, în contextul utilizării procesorului hibrid în aplicații DSP, viteza de operare este aceeași datorită structurii pipeline în care este organizat;

- operația de împărțire, efectuată chiar și singular, este de circa 3 ori mai rapidă decât una în virgulă flotantă;

- subunitatea logaritmice permite și implementarea cu ușurință a unor operații de extragere a rădăcinii pătrate, de ridicare la orice putere reală, sau de logaritmare în orice bază, care se pot efectua mult mai rapid decât în virgulă flotantă (de 5÷10 ori).

Evaluarea cea mai potrivită a performanțelor sistemului de calcul proiectat se face însă prin compararea cu ultimele realizări prezentate în literatura de specialitate în domeniul procesoarelor aritmetice hibride, care aparțin lui F. S. Lai și E. Wu. Față de aceste referințe, similare în performanțe, procesorul aritmetic hibrid de concepție proprie este de 1,6 ori mai rapid. Obținerea acestor performanțe a fost posibilă datorită unui șir de contribuții teoretice și aplicative, pe care autorul tezei le-a adus prin elaborarea și implementarea unor metode specifice de operare numerică.

Ca și contribuții teoretice au fost considerate și metodele originale aplicate în cazurile particulare de proiectare a unor blocuri funcționale concrete, dar care prezintă un grad de generalitate ce permite utilizarea lor și în alte aplicații de proiectare hardware. Contribuțiile aduse de autorul tezei pe parcursul acesteia vor fi menționate în continuare.

A. Contribuții teoretice:

1. Introducerea formatului LNS simplă precizie cu deplasament, ce a permis modificarea parțială a algoritmului descris și utilizat de F. S. Lai în [50] și [107] (v.

paragraful 2.2.1, relațiile 2.14-2.15 și explicațiile aferente). Partea întreagă a logaritmului unui operand are astfel inclusă valoarea 127, ceea ce asigură operanzi de intrare în blocul ALU al subunității logaritmice, întotdeauna pozitivi. Utilizarea formatului LNS cu deplasament a condus la economie de hard prin eliminarea unor sumatoare pentru exponenți și prin simplificarea blocului ALU și deci la creșterea vitezei de operare a acestuia.

2. Introducerea metodei “logaritmului mascat” ca metodă de creștere a vitezei de operare în cazul utilizării unui algoritm oarecare nerecursiv ce presupune calculul unei sume de produse de factori și de alți termeni distincți, prin considerarea tuturor termenilor ca și pseudo-produse parțiale și includerea lor într-un arbore Wallace (v. paragraful 2.2.3.1).

3. Îmbunătățirea algoritmului de interpolare liniară între valori memorate succesive ale unor funcții continue nemonotone, prin eliminarea operațiilor de scădere a unor termeni (prin memorarea directă a complementului de doi a scăzătorului) sau a unor produse de factori corespunzătoare porțiunilor coborâtoare ale funcțiilor (în urma translatării cu o locație pe aceste porțiuni a conținutului memoriilor ROM utilizate) (v. paragraful 2.2.3.1, relațiile 2.19, 2.20, fig. 2.16 și explicațiile aferente). În aplicația concretă, acest lucru a condus la creșterea nemijlocită a vitezei de operare a circuitelor de logaritmare și antilogaritmare cât și la obținerea unei structuri hard identice pentru cele două circuite, diferențele regăsindu-se în conținutul memoriilor ROM utilizate).

4. Demonstrarea posibilității de reducere a erorilor de conversie la logaritmare și antilogaritmare, prin adăugarea unor valori de corecție în anumite locații ale memoriilor ROM utilizate (v. paragraful 2.2.4, fig. 2.26 și explicațiile aferente).

5. Demonstrarea posibilității de reducere a dimensiunii structurilor numerice care implementează algoritmi de interpolare liniară pe baza calculării erorii de trunchiere cumulate în cazul eliminării celor mai puțin semnificative poziții de biți din arborele Wallace folosit (v. paragraful 2.2.4, explicațiile referitoare la fig 2.19).

6. Redefinirea compresorului 4:2 pe baza ieșirilor sale redundante, ceea ce a permis lărgirea familiei de circuite din care face parte (v. paragraful 2.2.3.2).

7. Corectarea formulelor date de W. Stallings în [80] care permit determinarea timpului maxim de generare a rezultatului la sumatoarele cu anticiparea transportului și cu selecția transportului (v. paragraful 1.1, relațiile 1.9, 1.10 și explicațiile aferente).

B. Contribuții aplicative:

1. Arhitectură (inclusiv organizare pe niveluri pipeline în două variante) de concepție proprie a subunității logaritmice (v. fig. 2.3, fig. 2.14 și fig. 2.25).

2. Determinarea dimensiunii memoriilor ROM utilizate la calculul logaritmului și antilogaritmului (v. paragraful 2.2.2.1).

3. Topologie de concepție proprie pentru circuitul sumator/scăzător, obținută prin schimbarea amplasamentului a două blocuri selectoare din structura clasică, ce

a condus la scurtarea căii critice de propagare prin circuit (v. fig. 2.9 și explicațiile aferente).

4. Modificarea structurii clasice a blocurilor componente ale sumatorului cu anticiparea transportului (CLA) cu structură modulară, ce a condus la creșterea vitezei de generare a rezultatului și la un număr mai mic de porți logice utilizate (v. fig. 2.10.a, fig. 2.10.b, fig. 2.10.c, fig. 2.11 și explicațiile aferente).

5. Conceperea unui program scris în limbaj C pentru depistarea celui mai defavorabil caz de propagare a transportului prin sumatorul de tip CLA (v. pag. 55, aliniatul 2 și Anexa 3).

6. Conceperea unui nou tip de sumator cu selecția transportului organizat pe două niveluri ierarhice, la care blocurile componente sunt sumatoare de tip CLA (v. paragraful 2.2.2.4, fig. 2.11, fig. 2.12 și explicațiile aferente)

7. Definirea căilor de selecție prioritare la sumatorul menționat anterior, ce a condus la o creștere suplimentară a vitezei acestuia (v. fig. 2.13 și explicațiile aferente).

8. Găsirea și explicarea celui mai defavorabil caz la propagarea transportului prin sumatorul menționat anterior (v. paragraful 4.1, pag. 97-98).

9. Construcția nivelului pipeline 2 al subunității logaritmice prin implementarea metodei logaritmului mascat și a metodei de interpolare liniară îmbunătățite (v. paragraful 2.2.3.1, fig. 2.17, fig. 2.18, fig. 2.19, fig. 2.20 și explicațiile aferente).

10. Conceperea blocului ALU al subunității logaritmice în cazul implementării metodei logaritmului mascat (v. paragraful 2.2.2.3, relațiile 2.25÷2.29, fig. 2.24 și explicațiile aferente).

11. Conceperea compresorului 4:2 L, mai rapid decât variantele descrise în literatura de specialitate (v. paragraful 2.2.3.2, relațiile 2.23-2.24, fig. 2.23 și explicațiile aferente).

12. Conceperea blocului logic ce controlează funcționarea ALU din subunitatea de calcul în virgulă flotantă (pentru toate combinațiile posibile din punct de vedere al semnelor operanzilor și al operației care se execută) și distribuția sa între nivelele pipeline 1 și 2, cu efect în creșterea vitezei de operare, respectiv în optimizarea din punct de vedere al suprafeței ocupate pe cip a structurii blocului ALU (paragraful 3.2, relațiile 3.1÷3.4, fig. 3.7 și explicațiile aferente).

13. Optimizarea circuitelor de aliniere a mantiselor, de adunare/scădere și de normalizare din subunitatea de calcul în virgulă flotantă în scopul reducerii suprafeței ocupate pe cip, în condițiile garantării sincronizării (funcționării cu aceeași perioadă de tact) cu subunitatea de calcul logaritmică (v. cap.3 în întregime).

14. Elaborarea diagramelor ce furnizează dimensiunea operanzilor, în cazul unui format nestandardizat, atunci când se realizează înmulțirea logaritmică prin citirea directă în memoria ROM, iar eroarea maximă de calcul este impusă (v. paragraful 1.3.1, relația 1.26, fig. 1.13, tab. 1.2 și explicațiile aferente).

15. Conceperea unui program Matlab pentru generarea valorilor ce trebuie memorate în ROM, în condițiile menținerii erorii de conversie la logaritmare și antilogaritmare sub $1,8 \times 10^{-7}$.

Toate blocurile funcționale ale subunității logaritmice și ale subunității de calcul în virgulă flotantă au fost proiectate cu ajutorul programului MicroSim până la nivel de poartă logică sau poartă de transmisie și apoi simulate, confirmându-se pe părți și pe ansamblu funcționarea corectă și în parametrii de viteză menționați.

Evaluarea timpului de procesare a informației s-a efectuat în condițiile utilizării tehnologiei CMOS $0,5\mu\text{m}$. Utilizarea însă a unei tehnologii submicronice mai avansate va permite operarea la o frecvență de clock a sistemului de calcul proiectat ce se apropie de 1 GHz ceea ce îl face competitiv cu cele mai recente realizări din domeniu.

Pe cadrul reprezentat de sistemul de calcul conceput în această lucrare poate fi adăugat hardware pentru realizarea unor funcții specifice, dedicate unui anumit context matematic oferit de aplicația concretă în care se dorește utilizarea procesorului aritmetic. Acest fapt deschide și perspectivele de dezvoltare ulterioară a temei abordate, susținute și de avansul tehnologic actual:

- conceperea unui mecanism de acumulare la nivelul ALU din subunitatea logaritmă pentru efectuarea eficientă a unor produse sau fracții multiple;

- conceperea unui procesor aritmetic puternic pentru aplicații DSP ce solicită volum mare de calcul în timp real și care presupune operarea paralelă în mai multe subunități logaritmice sau în mai multe circuite de logaritmare, pentru creșterea numărului de termeni ce se calculează simultan, respectiv a numărului de factori ai produselor calculate. În ambele situații utilizarea metodei logaritmului mascat va permite creșterea vitezei de operare și economie de hard prin utilizarea unui singur acumulator;

- concepția unui procesor aritmetic de uz general, caz în care implementarea operațiilor de extragere a rădăcinii pătrate, ridicare la putere și logaritmare presupune introducerea în ALU a unui circuit de deplasare precum și a unui sumator într-unul din cele două circuite de logaritmare, pentru extragerea deplasamentului din exponentul operandului. De asemenea va trebui adăugat hard suplimentar în ambele subunități ale procesorului aritmetic pentru introducerea mecanismelor de rotunjire a rezultatului în conformitate cu standardul IEEE754.

- reorganizarea subunității de calcul în virgulă flotantă pe doar două niveluri pipeline, prin secționarea sumatoarelor de tip CLA din ALU și aplicarea metodelor menționate în paragraful 1.2.3;

- în cazul unei posibilități de implementare practică pe cip, proiectarea circuitelor de distribuție a clock-ului pentru latch-urile dintre nivelurile pipeline și desriere VHDL pentru toată structura hard proiectată;

- concepția unui procesor aritmetic hibrid care să opereze în dublă precizie.

Bibliografie

- [1] Al-Twaijry H. A., Flynn M. J., *Technology Scaling Effects on Multipliers*, IEEE Transactions on Computers, Vol. 47, No. 11, Nov. 1998.
- [2] Arnold M. G., Bailey T., Cowles J., *Arithmetic Co-Transformations in the Real and Complex Logarithmic Number System*, IEEE Transactions on Computers, Vol. 47, No. 7, July 1998.
- [3] Arnold M. G., Bailey T., Cowles J., *Applying Features of IEEE754 to Sign/Logarithm Arithmetic*, IEEE Transactions on Computers Vol. 41, No.8, August 1992.
- [4] Arnold M. G., Bailey T., Cowles J., *Redundant Logarithmic Arithmetic*, IEEE Transactions on Computers Vol.39 NO.8, August 1990.
- [5] Arnold M., Bailey T., Cowles J., *Comments on "An Architecture for Addition and Subtraction of Long Word Length Numbers in the Logarithm Number System"*, IEEE Transactions on Computers, Vol. 41, No. 6, June 1992.
- [6] Bârsan R. M., *Fizica și tehnologia circuitelor MOS integrate pe scară mare*, Editura Academiei R.S.R., București, 1989.
- [7] Beaumont-Smith A., Burgess N., Lefrere S., Lim C.C., *Reduced Latency IEEE Floating-Point Standard Adder Architectures*, 14th IEEE Symposium on Computer Arithmetic, Adelaide, April 1999.
- [8] Benschneider B. J., Bowhill W. J., *A 50MHz Uniformly Pipelined 64b Floating-Point Arithmetic Processor*, IEEE International Solid-State Circuits Conference, 1990.
- [9] Brubaker T. A., Becker J. C., *Multiplication Using Logarithms Implemented with Read-Only Memory*, IEEE Transactions on Computers Vol. C-24 No.8 August 1975.
- [10] Bruguera J. D., Lang T., *Leading-One Prediction Scheme for Latency Improvement in Single Datapath Floating-Point Adders*, Proc. Intel Conf. Computer Design ICCD'98, Austin, Texas, 1998.
- [11] Bruguera J. D., Lang T., *Leading-One Prediction with Concurrent Position Correction*, IEEE Transactions on Computers, Vol. 48, No. 10, October 1999.
- [12] Chen. C. & all, *Pipelined Computation of Very Large Word Length LNS Addition/Subtraction with Polynomial Hardware Cost*, IEEE Transactions on Computers, Vol. 49, No. 7, July 2000.
- [13] Cheng F., Unger S. H., Theobald M., *Self-Timed Carry-Lookahead Adders*, IEEE Transactions on Computers, Vol. 49, No. 7, July 2000.
- [14] Ciminiera L., Montuschi P., *Carry-Save Multiplication Schemes Without Final Addition*, IEEE Transaction on Computers, Vol.45, No.9, September 1996.
- [15] Ciugudean M., *Algoritmi pentru operații cu logaritmi binari în unități aritmetice cu virgulă flotantă*, Simpozionul Național de Teoria Sistemelor, Universitatea Craiova, 1986.

[16] Ciugudean M., *Algoritmi și structuri logice celulare de generare a logaritmilor și antilogaritmilor binari pentru dispozitive aritmetice cu virgulă mobilă*, Teză doctorat, I.P. T.V., Timișoara, 1976.

[17] Ciugudean M., *Structures logiques cellulaires pour la génération des logarithmes et antilogarithmes binaires*, Buletinul I.P.T.V. Timișoara, Tom 23(37), Fascicola 1, 1978.

[18] Ciugudean M., *Unitate aritmetică în virgulă flotantă logaritmică*, Simpozionul Național de Teoria Sistemelor, Universitatea Craiova, 1986.

[19] Clouser J. & all, *A 600-MHz Superscalar Floating-Point Processor*, IEEE Journal of Solid-State Circuits, Vol. 34, No. 7, July 1999.

[20] Coleman J., Chester E., Softley C., Kadlec J., *Corrections to "Arithmetic on the European Logarithmic Microprocessor"*, IEEE Transactions on Computers, Vol. 49, No. 10, Oct. 2000.

[21] Coleman J., Chester E., Softley C., Kadlec J., *Arithmetic on the European Logarithmic Microprocessor*, IEEE Transactions on Computers, Vol. 49, No. 7, July 2000.

[22] Das D. & all, *Implementation of Four Common Functions on an LNS Co-Processor*, IEEE Transactions on Computers, Vol. 44, No. 1, January 1995.

[23] Ercegovic M. D. & all, *Reciprocation, Square-Root and Some Elementary Functions Using Small Multipliers*, IEEE Transactions on Computers, Vol. 49, No. 7, July 2000.

[24] Fenn S. T. J., Benaissa M., Taylor D., *Multiplication and Division Over Dual Basis*, IEEE Transaction on Computers, Vol.45, No.3, March 1996.

[25] Fiore P. D., *Parallel Multiplication Using Fast Sorting Networks*, IEEE Transactions on Computers, Vol. 48, No.6, June 1999.

[26] Fujii H. & all, *A Floating-Point Cell Library and a 100-MFLOPS Image Signal Processor*, IEEE Journal of Solide-State Circuits, Vol.27, No.7. July 1992.

[27] Goto J. & all, *250-MHz BiCMOS Super-High-Speed Video Signal Processor (S-VSP) ULSI*, IEEE Journal of Solid-State Circuits, Vol. 26, No. 12, December 1991.

[28] Guyot A., Muller J., *A Way to Build Efficient Carry-Skip Adders*, IEEE Transactions on Computers, Vol. C-36, No. 10, October 1987.

[29] Hanyu T., Kameyama M., *A 200 MHz Pipelined Multiplier Using 1.5 V-Supply Multiple-Valued MOS Current-Mode Circuits with Dual-Rail Source-Coupled Logic*, IEEE Journal of Solid-State Circuits, Vol. 30, No. 11, November 1995.

[30] Hiasat A. A., *New Efficient Structure for a Modular Multiplier for RNS*, IEEE Transactions on Computers, Vol. 49, No. 2, Feb. 2000.

[31] Hummel R., *Processeur et Coprocesseur*, Editeur DUNOD, 1992.

[32] Ide N. & all, *2.44-GFLOPS 300-MHz Floating-Point Vector-Processing Unit for High-Performance 3-D Graphics Computing*, IEEE Journal of Solid-State Circuits, Vol. 35, No. 7, July 2000.

Bibliografie

- [33] Ide N., Kondo H.F., *A 320-MFLOPS CMOS Floating-Point Processing for Superscalar Processors*, IEEE Journal of Solid-State Circuits, Vol.28, No.3, March 1993.
- [34] Inoue T. & all, *A 2000-MOPS Embedded RISC Processor with a Rambus DRAM Controller*, IEEE Journal of Solid-State Circuits, Vol. 34, No. 7, July 1999.
- [35] Itoh N. & all, *A 600-MHz 54×54-bit Multiplier with Rectangular-Styled Wallace Tree*, IEEE Journal of Solid-State Circuits, Vol. 36, No. 2, February 2001.
- [36] Jessani R. M., Putrino M., *Comparison of Single- and Dual-Pass Multiply-Add Fused Floating-Point Units*, IEEE Transactions on Computers, Vol. 47, No. 9, Sept. 1998.
- [37] Jurca L. A., *Conceptia unui procesor aritmetic cu virgulă flotantă logaritmic*, Referat Nr.2, Timișoara, 1997.
- [38] Jurca L. A., Drăgoi B., *Using 4:2 L Compressor in Partial Product Reduction Tree for Fast Binary Multipliers*, Buletinul UPT, Tom 46(60), Fasc. 1, 2001.
- [39] Jurca L. A., *Floating-Point Unit for Fast Performing Addition and Subtraction in Single Precision Format*, Buletinul UPT, Tom 45(59), Fasc. 2, 2000.
- [40] Jurca L. A., *Logarithmic Unit for Fast Performing of Multiplication and Division Operations in Single Precision Format*, Buletinul UPT, Tom 45(59), Fasc. 2, 2000.
- [41] Jurca L. A., *Method and Circuit for the Fast Calculation of the Binary Logarithm and Antilogarithm in Single Precision Format*, Buletinul UPT, Tom 45(59), Fasc.1, 2000.
- [42] Jurca L. A., *Some Considerations Regarding the Design of a Hybrid Logarithmic, Floating-Point Mathematical Processor*, Buletinul UPT, Tom 45(59), Fasc. 1, 2000.
- [43] Jurca L. A., *Stadiul actual al procesoarelor aritmetice cu virgulă flotantă*, Referat Nr.1, Timișoara, 1997.
- [44] Jurca L. A., *Using Masked Logarithm Method for Fast Multiplication and Division Logarithmic Unit Design*, Buletinul UPT, Tom 46(60), Fasc. 1, 2001.
- [45] Kadlec J., Softley C. & all, *32-bit Logarithmic ALU for Handel C 2.1 and Celoxica DK1 (53 MHz for XCV2000E-6 based RC1000 board)*, Celoxica User Conference, Stratford, UK, April 2001.
- [46] Kantabutra V., *Accelerated Two-Level Carry-Skip Adders – A Type of Very Fast Adders*, IEEE Transactions on Computers, Vol. 42, No. 11, November 1993.
- [47] Kantabutra V., *Designing Optimum One-Level Carry-Skip Adders*, IEEE Transactions on Computers, Vol. 42, No. 6, June 1993.
- [48] Kantabutra V., *On Hardware for Computing Exponential and Trigonometric Functions*, IEEE Transactions on Computers, Vol.45, No.3, March 1996.

Bibliografie

[49] Karagianni K. & all, *Operation-Saving VLSI Architectures for 3D Geometrical Transformation*, IEEE Transactions on Computers, Vol. 50, No. 6, June 2001.

[50] Lai F., *A 10-ns Hybrid Number System Data Execution Unit for Digital Signal Processing Systems*, IEEE Journal of Solid-State Circuits, Vol. 26, No. 4, Apr. 1991.

[51] Lai F., Wu C.F.E., *A Hybrid Number System Processor with Geometric and Complex Arithmetic Capabilities*, IEEE Transactions on Computers Vol. 40 No.8 Aug.1991.

[52] Lang J. H. & all, *Integrated-Circuit Logarithmic Arithmetic Units*, IEEE Trans. Comp Vol C-34 No.5, August 1985.

[53] Lang T., Montuschi P., *Boosting Very-High Radix Division with Prescaling and Selection by Rounding*, IEEE Transactions on Computers, Vol. 50, No. 1, January 2001.

[54] Lang. T., Montuschi P., *Verry High Radix Square Root with Prescaling and Rounding and a combined Division/Square Root Unit*, IEEE Transactions on Computers, Vol. 48, No. 8, August 1999.

[55] Larsson-Edefors P., *A 965-Mb/s 1.0- μ m Standard CMOS Twin-Pipe Serial/Parallel Multiplier*, IEEE Journal of Solid-State Circuits, Vol. 31, No. 2, February 1996.

[56] Lev L. A. & all, *A 64-b Processor with Multimedia Support*, IEEE Journal of Solid-State Circuits, Vol. 30, No. 11, November 1995.

[57] Lewis D. M., *114 MFLOPS LNS Arithmetic Unit for DSP Applications*, IEEE Journal of Solide-State Circuits, vol.30, No.12, December 1995.

[58] Lewis D. M., *An Architecture for Addition and Subtraction of Long Word Length Numbers in the Logarithm Number System*, IEEE Transactions on Computers Vol. 39, No.11, Nov.1990.

[59] Lewis D. M., *Interleaved Memory Function Interpolators with Application to Accurate LNS Arithmetic Unit*, IEEE Transactions on Computers Vol. 43, No.8 August 1994.

[60] Lineback J. R., *Single Chip Processors*, Electronics, March 1988.

[61] Lo H., Aoki Y., *Generation of a Precise Binary Logarithm with Difference Grouping Programmable Logic Array*, IEEE Transactions on Computers Vol. C-34 No.8, August 1985.

[62] Lu F., Samuelli H., *A 200-MHz CMOS Pipelined Multiplier-Accumulator Using a Quasi-Domino Dinamic Full-Adder Cell Design*, IEEE journal of Solid-State Circuits, Vol.28, No.2, February 1993.

[63] Maenner R., *A Fast Integer Binary Logarithm of Large Arguments* IEEE Micro, December 1987.

[64] Minami T. & all, *A 300-MOPS Video Signal Processor with a Parallel Architecture*, IEEE Journal of Solid-State Circuits, Vol. 26, No. 12, December 1991.

Bibliografie

[65] Montoye R. K., Cook P. W., Hokenek E., Havvreluk R. P., *An 18ns 56-Bit Multiply-Adder Circuit*, IEEE International Solid-State Circuits Conference, 1990.

[66] Montoye R. K., Hokenek E., Runion S. L.: *Design of the IBM RISC Sistem/6000 Floating-Point Execution Unit*, IBM J. Res. Develop., Vol.34 No.1, January 1990.

[67] Mori J., Nagamatsu M. *A 10-ns 54×54-b Parallel Structured Full Array Multiplier with 0.5- μ m CMOS Tehnology*, IEEE Journal of Solid-State Circuits, Vol.26, No.4, April 1991.

[68] Nakayama T., Kojima S., *An 80b 6.7MFLOPS Floating-Point Processor with Vector / Matrix Instructions*, IEEE International Solid-State Circuits Conference, 1989.

[69] Oberman S. F., Flynn M. J., *Division Algorithms and Implementations*, IEEE Transactions on Computers, Vol. 46, No. 8, August 1997.

[70] Ohkubo N. & all, *A 4.4 ns CMOS 54×54 b Multiplier Using Pass-Transistor Multiplexer*, IEEE Journal of Solid-State Circuits, Vol. 30, No. 3, March 1995.

[71] Okamoto F. & all, *A 200-MFLOPS 100-MHz 64-b BiCMOS Vector-Pipelined Processor (VPP) ULSI*, IEEE Journal of Solid-State Circuits, Vol. 26, No. 12, December 1991.

[72] Oklobdzija V. G., Villeger D., Liu S. S., *A Method for Speed Optimized Partial Product Reduction and Generation of Fast Parallel Multipiliers Using an Algorithmic Approach*, IEEE Transaction on Computers, Vol.45, No.3, March 1996.

[73] Petterson D., Hennessy J. L., *Computer Architecture. A Quantitative Approach*, Morgan Press, 1996.

[74] Phatak D. S., Goff T., Koren I., *Constant-Time Addition and Simultaneous Format Conversion Based on Redundant Binari Representations*, IEEE Transactions on Computers, Vol.50, No. 11, Nov. 2001.

[75] Phatak D. S., Goff T., Koren I., *Redundancy Management in Arithmetic Processing via Redundant Binary Representations*, Asilomar '99.

[76] Schulte M.J., Swartzlander E.E., *Hardware Design for Exactly Rounded Elementary Functions*, IEEE Transactions on Computers Vol.43, No.8, Aug.1994.

[77] Seidel P.M., Even G., *An IEEE Floating-Point Adder Design Optimized For Speed*, 15th IEEE Symposium on Computer Arithmetic, Adelaide, April 2000.

[78] Senthinathan R. & all, *A 650-MHz, IA-32 Microprocessor with Enhanced Data Streaming for Graphics and Video*, IEEE Journal of Solid-State Circuits, Vol. 34, No. 11, November 1999.

[79] Srinivas H. R., Parhi K. K., *A fast VLSI Adder Architecture*, IEEE Journal of Solide-State Circuits, Vol.27, No.5, May 1992.

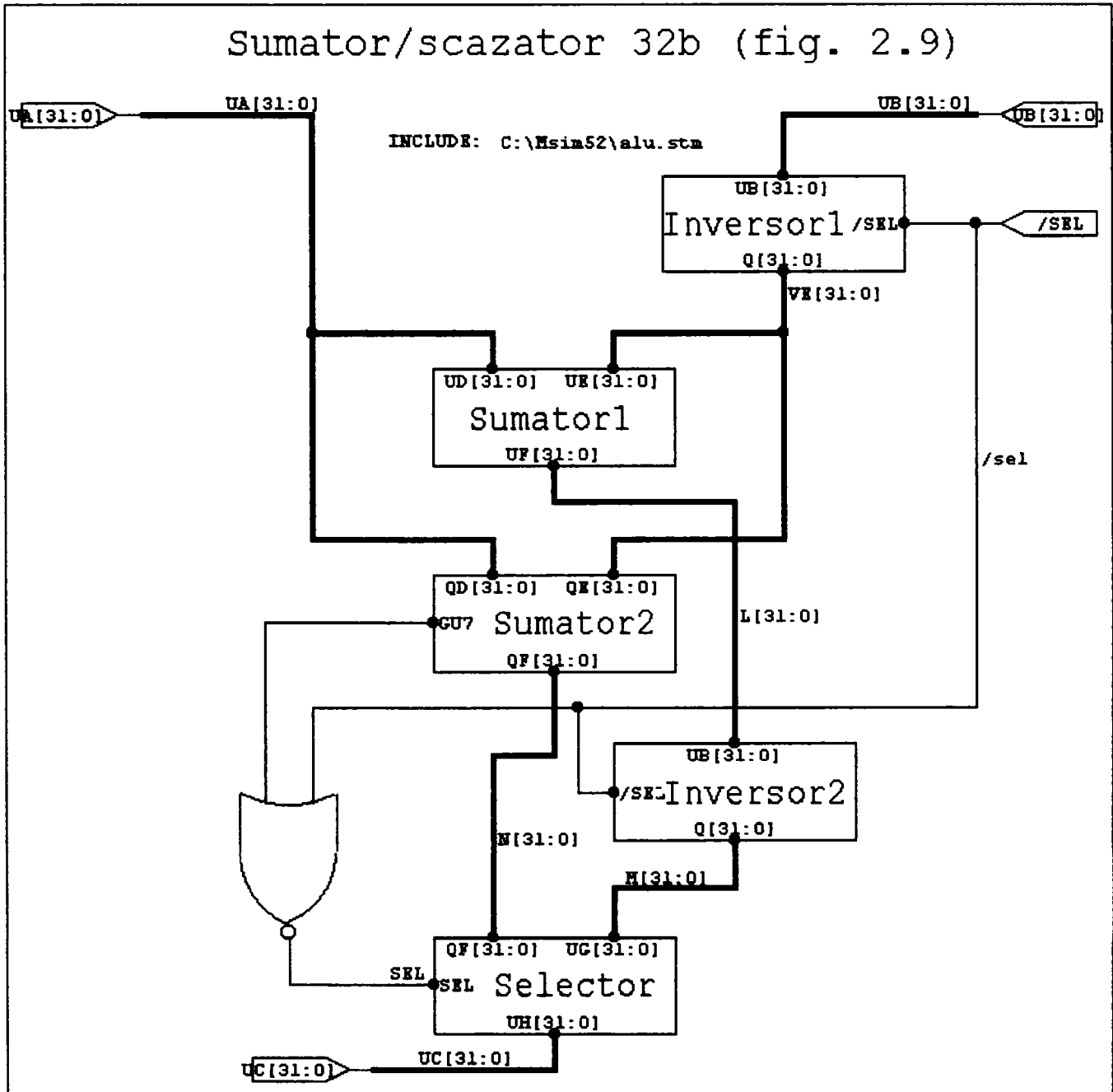
[80] Stallings W., *Computer Architecture and Organization*, Prentice Hall Inc, 1996.

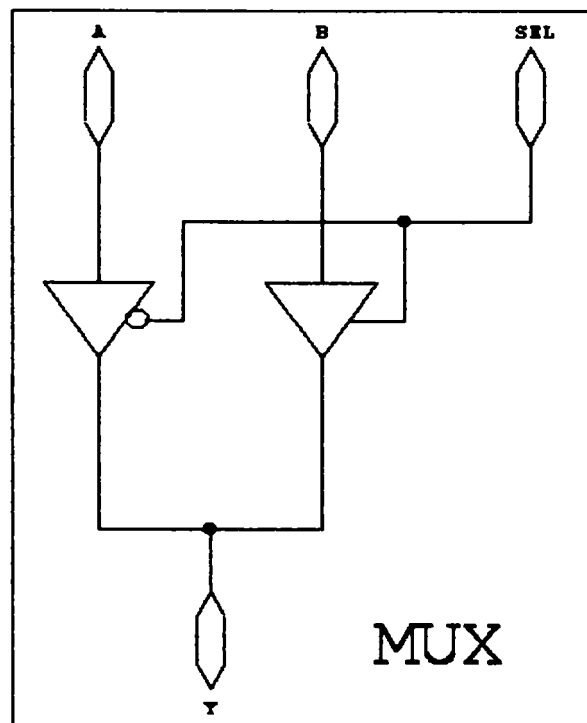
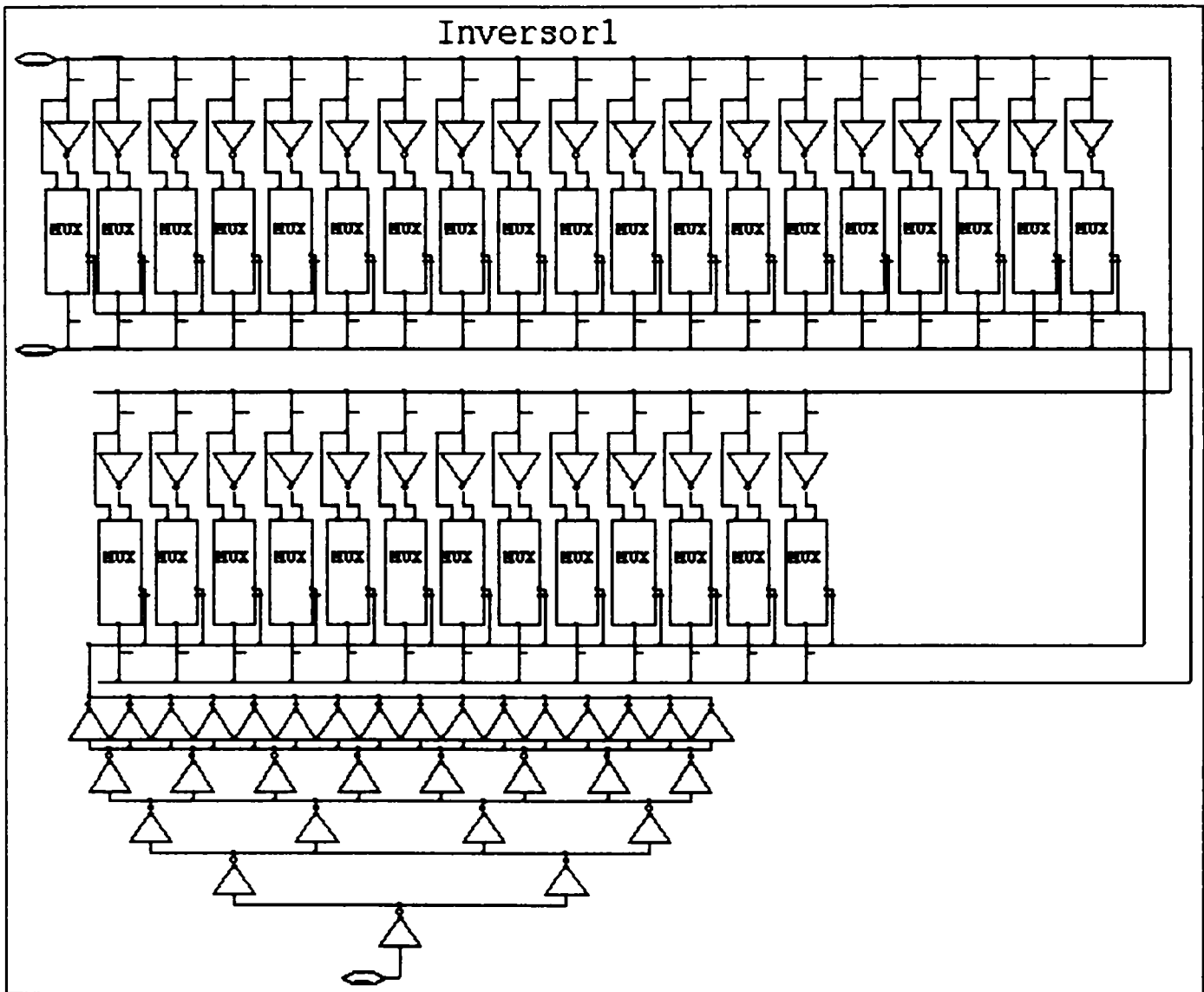
Bibliografie

- [81] Stelling P. F., Oklobdzija V. G. & all, *Optimal Circuits for Parallel Multipliers*, IEEE Transactions on Computers, Vol. 47, No. 3, March 1998.
- [82] Stouraitis T., Taylor F. J., *Floating-Point to Logarithmic Encoder Error Analysis*, IEEE Transaction on Computers, Vol.37, No.7, July 1988.
- [83] Strugaru C., Popa M., *Microprocesoare pe 16 biți*, Editura TM 1992.
- [84] Sunar B., Koc C., *An Efficient Optimal Normal Basis Type II Multiplier*, IEEE Transactions on Computers, Vol. 50, No.1, January 2001.
- [85] Suzuki H. & all, *Leading-Zero Anticipatory Logic for High-Speed Floating Point Addition*, IEEE Journal of Solid-State Circuits, Vol. 31, No. 8, Aug. 1996.
- [86] Suzuoki M., *A Microprocessor with a 128-Bit CPU, Ten Floating-Point MAC's, Four Floating-Point Dividers, and an MPEG-2 Decoder*, IEEE Journal of Solid-State Circuits, Vol. 34, No. 11, November 1999.
- [87] Swartzlander E. E., Alexopoulos A. G., *The Sign/Logarithm Number System*, IEEE Transactions on Computers December 1975.
- [88] Takagi N., *Powering by a Table Look-Up and a Multiplication with Operand Modification*, IEEE Transactions on Computers, Vol. 47, No. 11, Nov. 1998
- [89] Taylor F. J. & all, *A 20 bit Logarithmic Number System Processor*, IEEE Transactions on Computers Vol.37, No.2, Feb.1988.
- [90] Taylor F. J. *A Hibrid Floating - Point Logarithmic Number System Processor*, IEEE Trans. Circuits & Systems, Vol. CAS-32, January 1985.
- [91] Timmermann D., Rix B., Hahn H., Hosticka J., *A CMOS Floating-Point Vector-Arithmetic Unit*, IEEE Journal of Solid-State Circuits, Vol.29, No.5, May 1994.
- [92] Toacșe G., *Introducere în microprocesoare*, Editura Științifică și Enciclopedică, București, 1986.
- [93] Toacșe G., Nicula D., *Electronică Digitală*, Editura Teora, București, 1996.
- [94] Toma C. I., *Dispozitiv de înmulțire 4×2 integrat*, Buletinul I.P.T.V. Timișoara, Tom 21(35), Fascicola 1, 1976.
- [95] Um J., Kim T., *An Optimal Allocation of CSA in Arithmetic Circuits*, IEEE Transactions on Computers, Vol. 50, No. 3, March 2001.
- [96] Vinnakota B. *Implementing Multiplication with Split Read-Only Memory*, IEEE Transaction on Computers, Vol.44, No.11, November 1995.
- [97] Wang Z., Jullien H. A., Miller W. C., *A New Design Tehnique for Column Compression Multipliers*, IEEE Transactions on Computers, Vol.44, No.8, August 1995.
- [98] Weste N., Eshraghian K., *Principles of CMOS VLSI Design: A Systems Perspective*, Addison-Wesley Publishing Company, 1985.
- [99] Williams T., Patkar N., Shen G., *SPARC64 A64-b 64-Active-Instructions Out-of-Order-Execution MCM Processor*, IEEE Journal of Solid-State Circuits, Vol. 30, No. 11, November 1995.

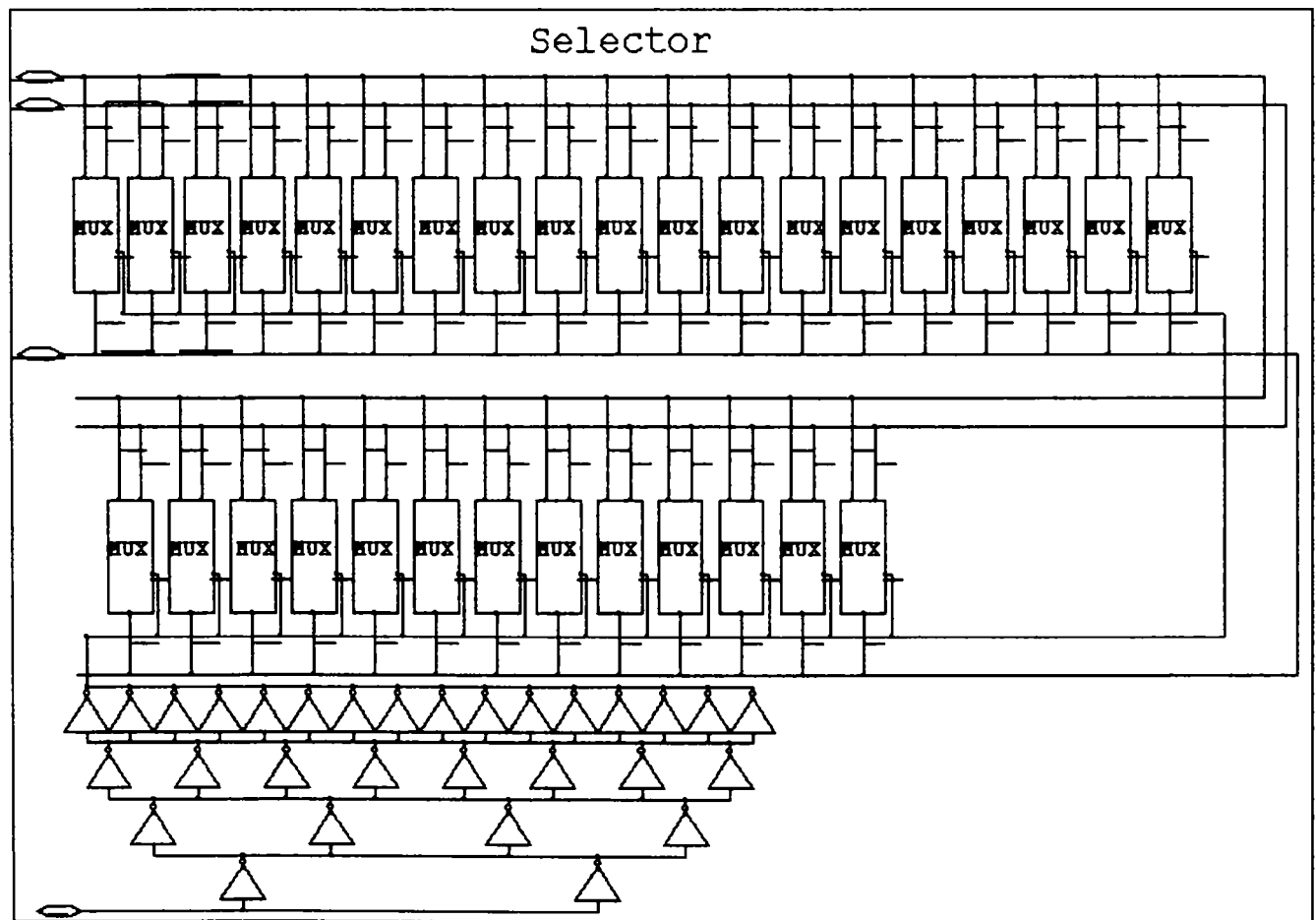
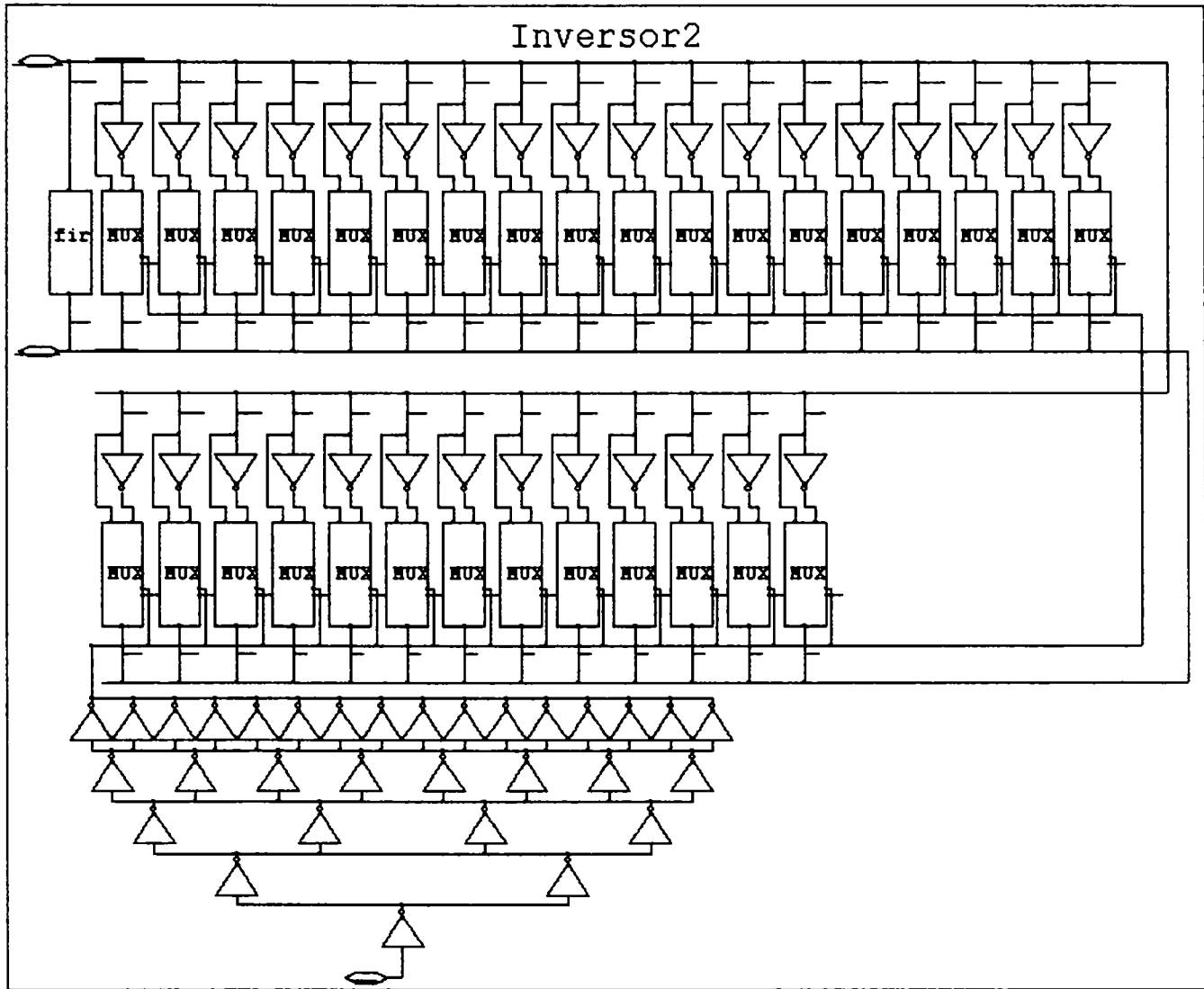
Bibliografie

- [100] Woods J. V. & all, *Amulet 1: An Asynchronous ARM Microprocessor*, IEEE Transactions on Computers, Vol. 46, No. 4, April 97.
- [101] Yan H. & all, *A Low-Power 16×16-b Parallel Multiplier Utilizing Pass-Transistor Logic*, IEEE Journal of Solid-State Circuits, Vol. 34, No. 10, October 1999.
- [102] Yang C., Sano B., Lebek A. R., *Exploiting Parallelism in Geometry Processing with General Purpose Processors and FP SIMD Instructions*, IEEE Transactions on Computers, Vol. 49, No. 9, Sept. 2000.
- [103] Yeh W., Jen C., *High-Speed Booth Encoded Parallel Multiplier Design*, IEEE Transactions on Computers, Vol. 49, No. 7, July 2000.
- [104] Yu L. K., Lewis D. M., *A 30-b Integrated Logarithmic Number System Processor*, IEEE Journal of Solid-State Circuits, Vol.26, No.10, October 1991.
- [105] Zhunang N., Wu H., *A new Design of the CMOS Full Adder*, IEEE Journal of Solide-State Circuits, Vol. 27, No.5, May 1992.
- [106] Villeger D., *Fast Parallel Multipliers*, Final Report, Ecole Superieure d'Ingenieurs en Electrotechnique et Electronique, Noisy-le-Grand, May 1993.
- [107] Lai F., *The Efficient Implementation and Analysis of a Hybrid Number System Processor*, IEEE Transactions on Circuits and Systems, Vol. 46, No. 6 ICSPE5, June 1993.
- [108] Butas J. & all, *Asynchronous Cross-Pipelined Multiplier*, IEEE Journal of Solid-State Circuits, Vol. 36, No. 8, August 2001.

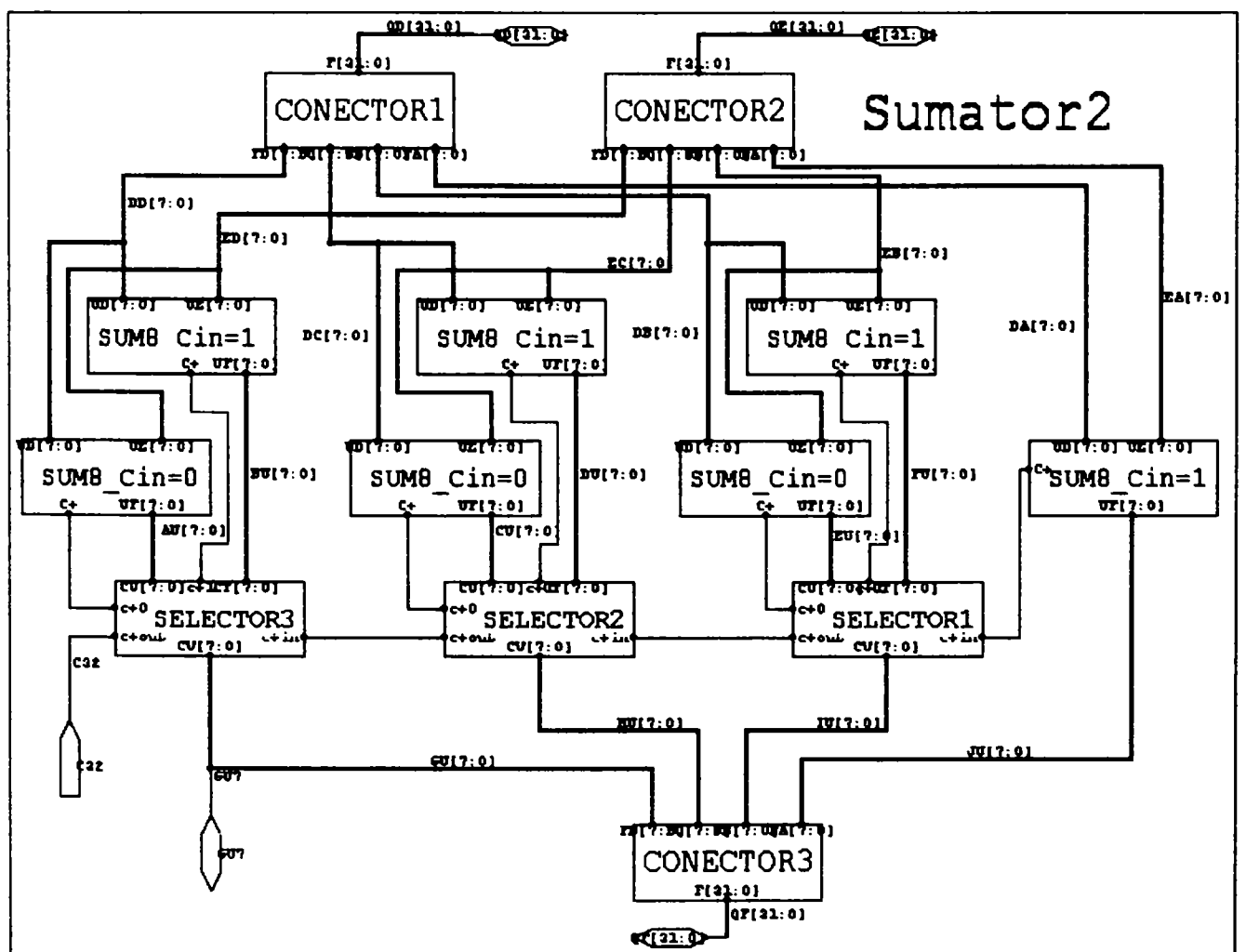
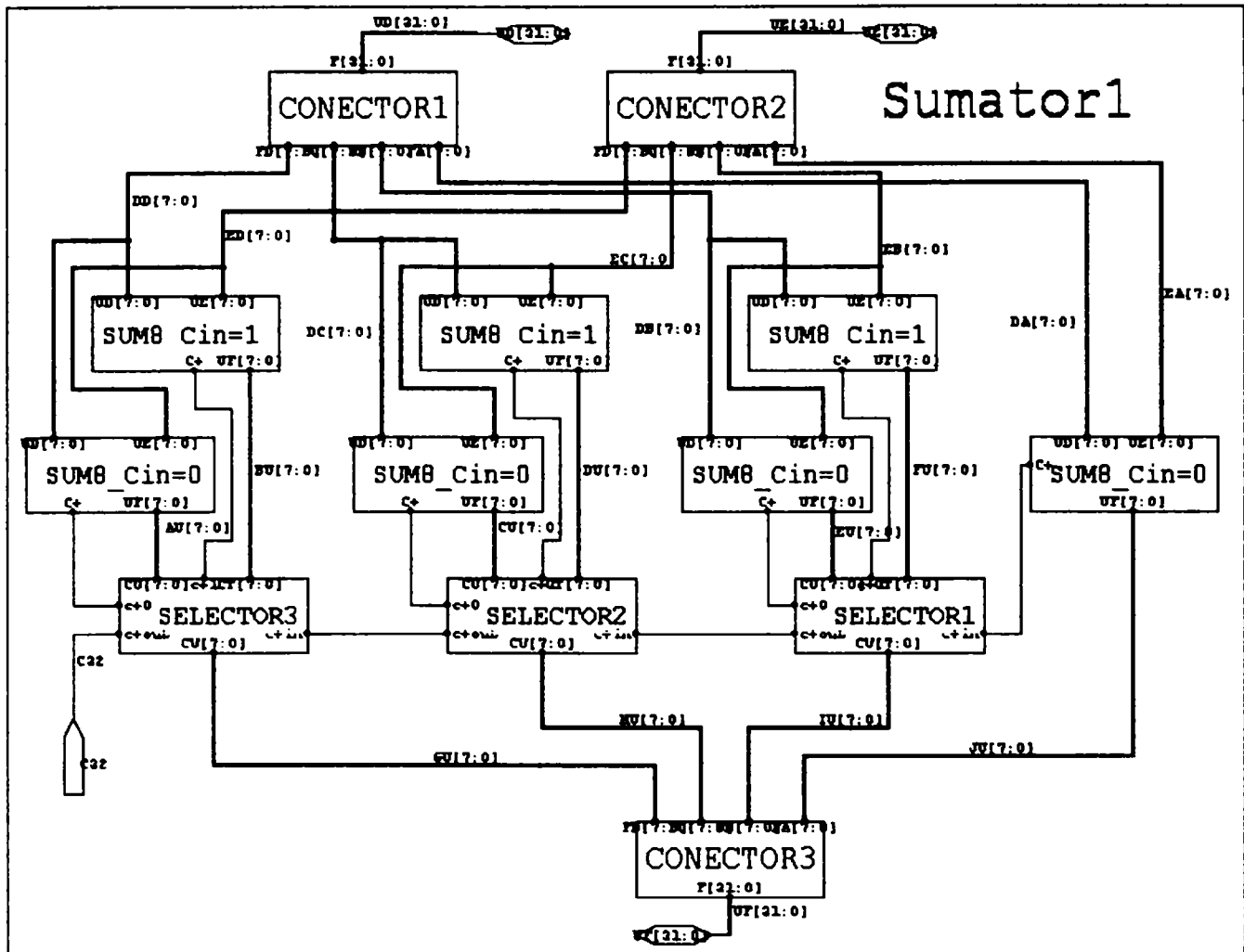


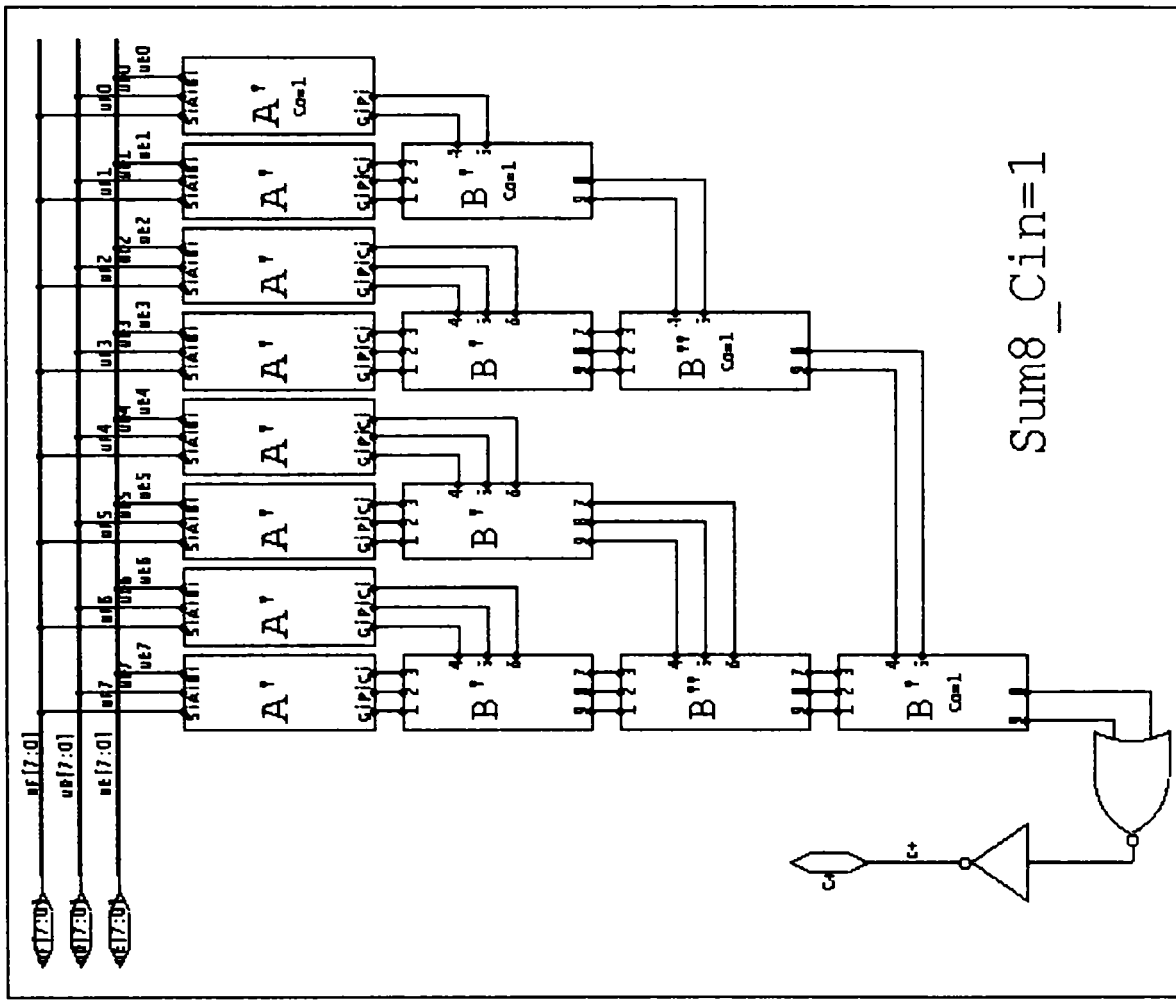
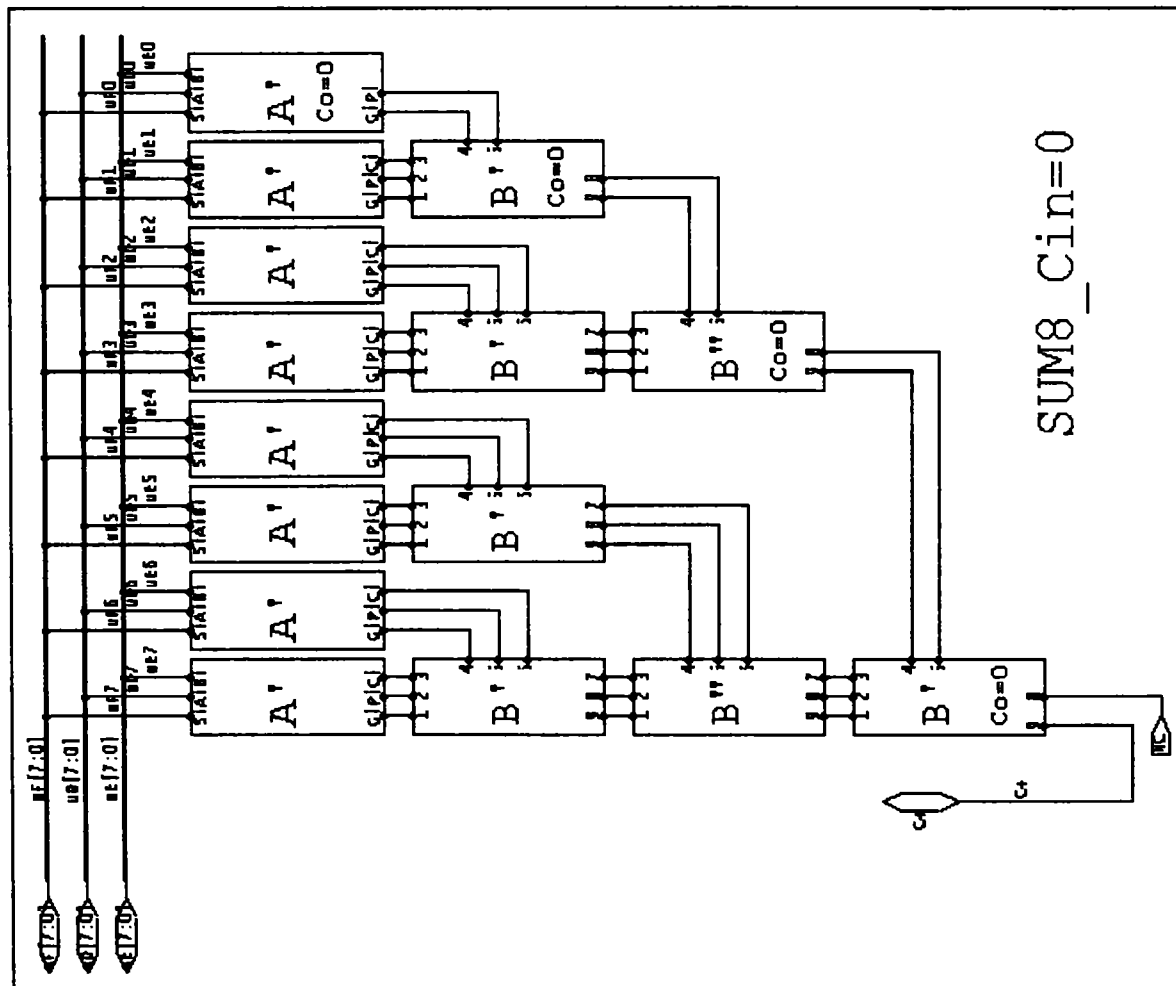


ANEXA 1



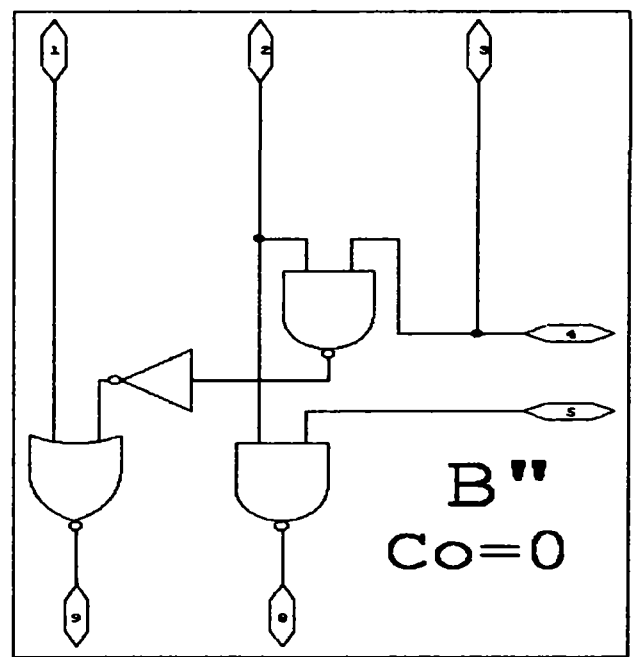
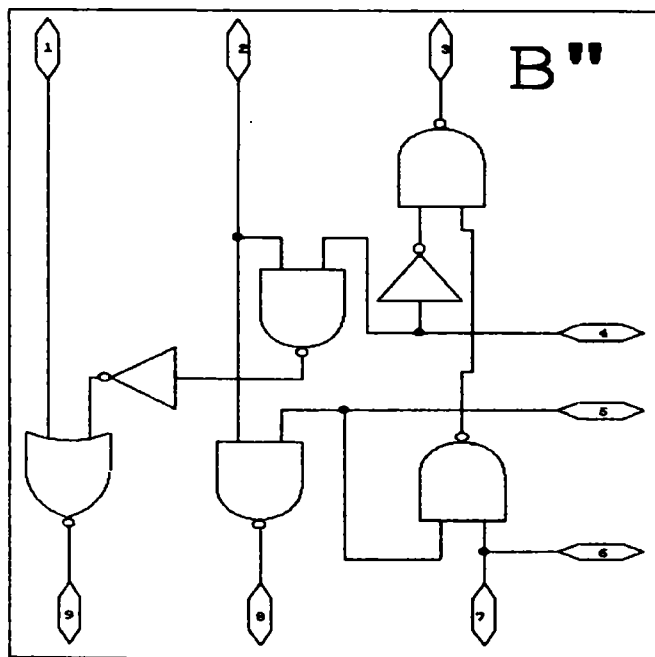
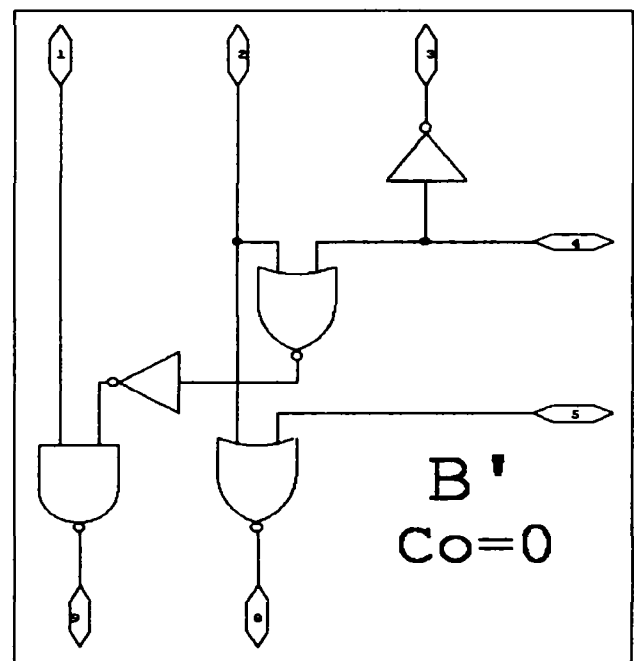
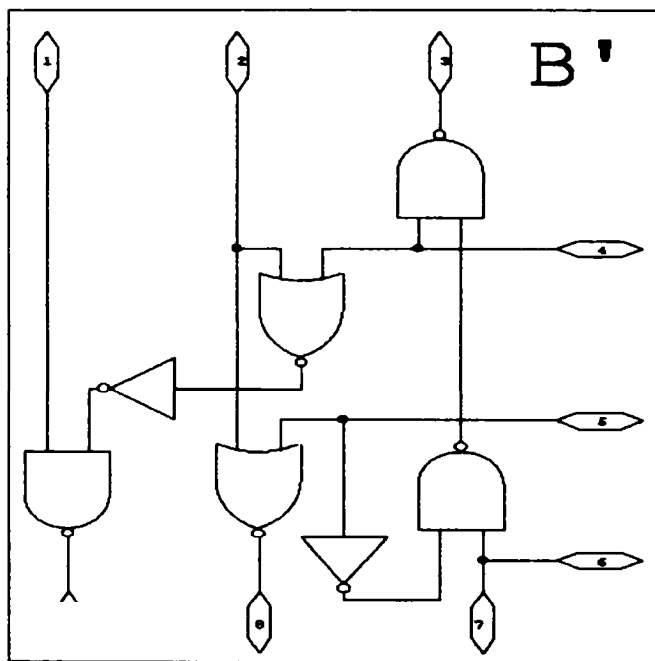
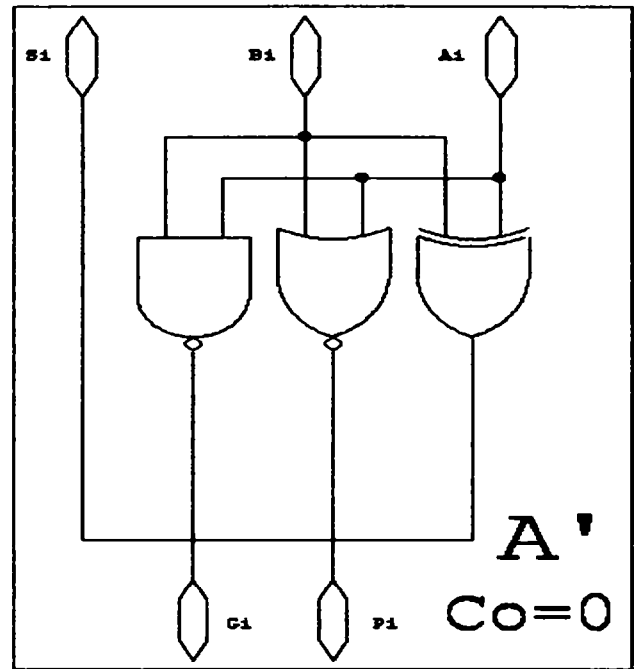
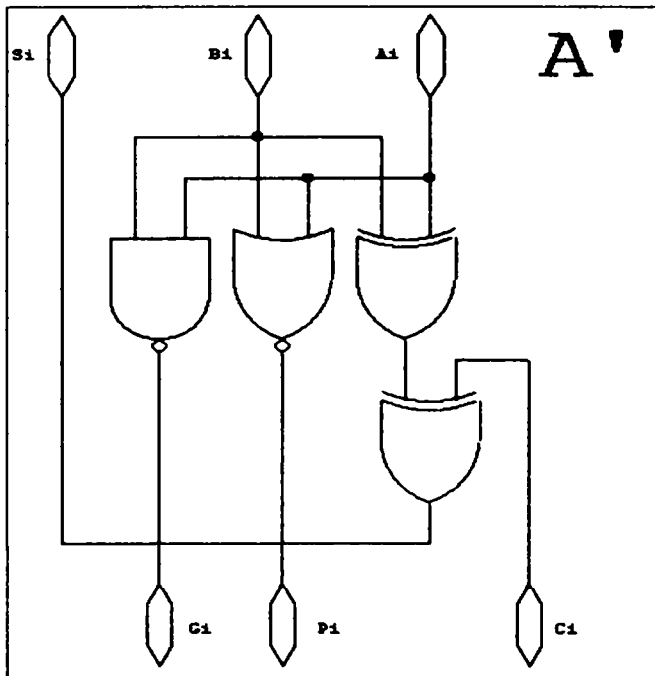
ANEXA 1



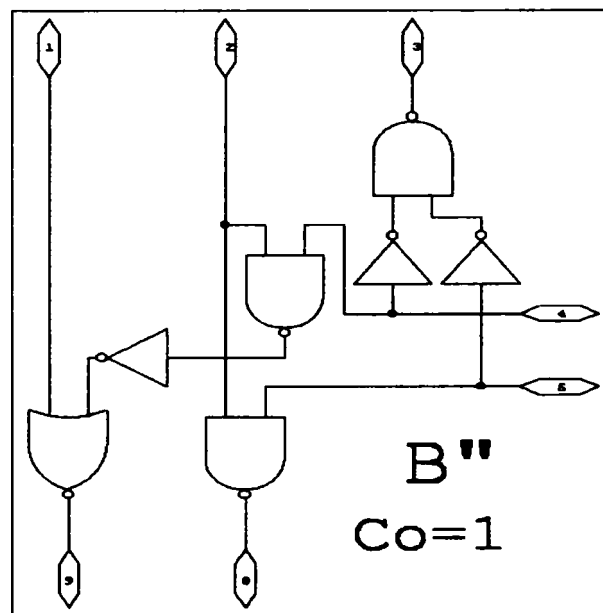
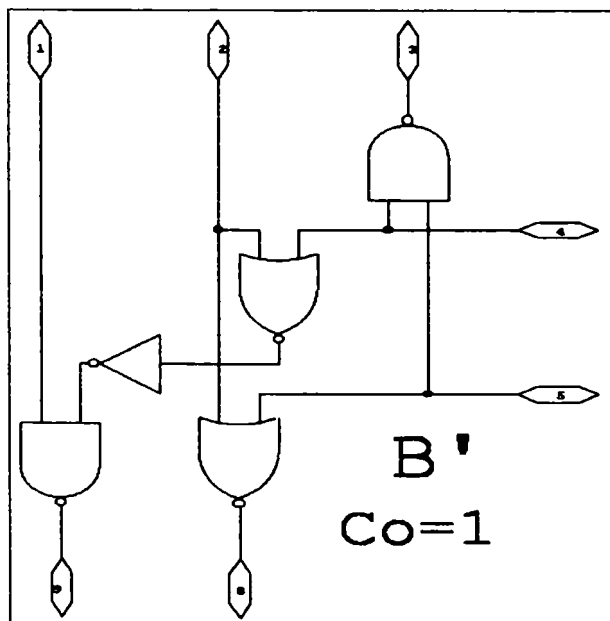
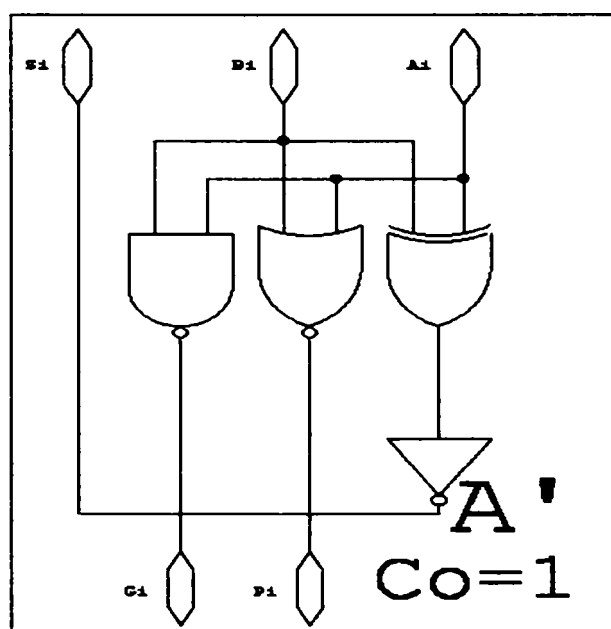
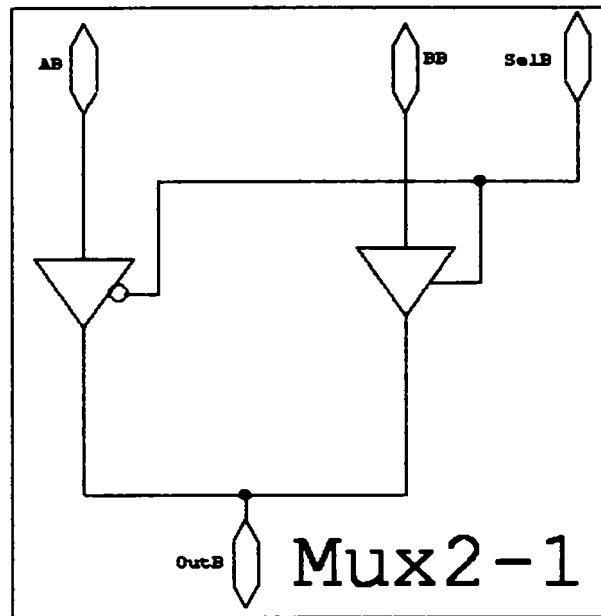
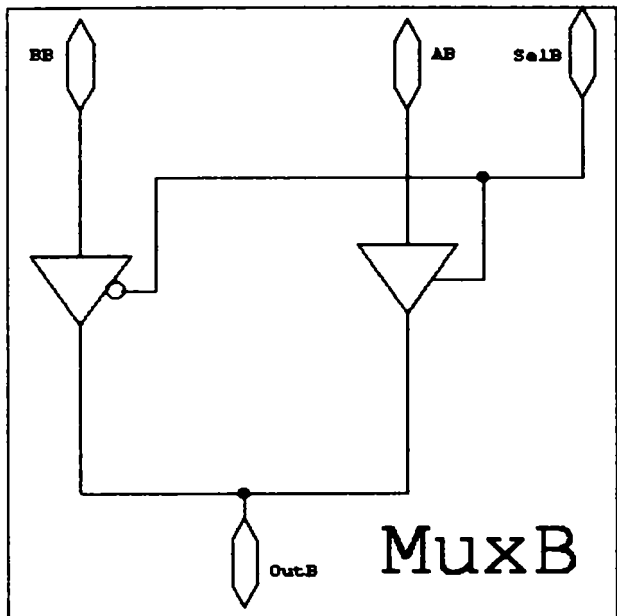


A1.5

ANEXA 1



ANEXA 1



ANEXA 2

```
#include <iostream.h>
#include <stdio.h>
#include <conio.h>
#include <stdlib.h>
#define marime 8
typedef struct
{
    char g,
        p,
        s,
        c,
        mem1,
        mem2,
        mem3;
} stare;
long double rang (int nb,int exp);
void genereaza_sir(long double a,char *ch);
char or(char a,char b);
char xor(char a,char b);
char and(char a,char b);
void init_stare(stare *gp1,stare *gp2,stare *gp3,stare *gp4);
int compara(stare stare_1,stare stare_2);
stare element_a(char x,char y, char c,stare a);
stare element_b(stare *gp, stare b,char c);
int carry_lookahead (long int _a_,long int _b_,long int _c_);

long double rang (int nb,int exp)
{
    long double p;
    int i;
    if(exp==0) // caz rang 0
    {
        p=1;
    }
    p=nb; // caz rang 1
    if(exp>1) // caz rang > 1
    {
        for(i=1;i<exp;i++)
        {
            p=nb*p;
        }
    }
    return(p); // trimita n^exp
}

void genereaza_sir(long double a,char *ch)
{
    int i=0;
    char x[marime+1];
    long double max;
    max=rang(2,(marime-1-i)); // calcul rang bit
    for(i=0;i<marime;i++)
    {
        if(0>a-max) // testeaza bitul i
            ch[i]='0';
        else
        {
            ch[i]='1';
            a=a-max;
        }
    }
}
```


ANEXA 2

```
    max=max/2;
}
ch[marime]='\0';           // plaseaza caracter de sfarsit de sir
}
char or(char a,char b)     // SAU
{
if((a==b)&(a=='0'))
{
return('0');
}
else
{
return('1');
}
}
char xor(char a,char b)   // SAU-exclusiv
{
if(a==b)
{
return('0');
}
else
{
return('1');
}
}
char and(char a,char b)   // SI
{
if((a==b)&(a=='1'))
{
return('1');
}
else
{
return('0');
}
}
void init_stare(stare *gp1,stare *gp2,stare *gp3,stare *gp4)
{
int i;
gp4->g='0';               // initializeaza fiecare camp al structurii
gp4->p='0';
gp4->mem1='0';
gp4->mem2='0';
gp4->mem3='0';
gp4->s='0';
gp4->c='0';
for(i=0;i<8;i++)        //initializeaza pe zero ceilalti g si p
{
gp1[i]=*gp4;
while(i<4)
{
gp2[i]=*gp4;
while(i<2)
{
gp3[i]=*gp4;
break;
}
}
break;
}
}
```

ANEXA 2

```

}
}

int compara (stare stare_1, stare stare_2)
{
    // compara fiecare camp al structurii
    int test=0;
    if(stare_1.g!=stare_2.g)test=1;    // in caz de diferenta pune test pe 1
    if(stare_1.p!=stare_2.p)test=1;
    if(stare_1.mem1!=stare_2.mem1)test=1;
    if(stare_1.mem2!=stare_2.mem2)test=1;
    if(stare_1.mem3!=stare_2.mem3)test=1;
    if(stare_1.c!=stare_2.c)test=1;
    if(stare_1.s!=stare_2.s)test=1;
    return(test);    // trimite test
}

stare element_a(char x,char y, char c, stare a)
{
    stare A;
    A.g=and(x,y);    // x&y=g
    A.p=or(x,y);    // p
    A.s=a.mem3;    //s    // + 1 timp de propagare=s
    A.mem3=xor(a.mem2,x);    // (y^c)^x
    A.mem2=a.mem1;    // + 1 timp de propagare
    A.mem1=xor(y,c);    // y^c
    A.c='\0';    // neutilizat
    return (A);    // trimite la
}

stare element_b(stare *gp, stare b ,char c)
{
    stare B;
    B.g=or(gp[0].g,b.mem1);    // g1|(p1&g2)=G
    B.mem1=and(gp[0].p,gp[1].g);    // p1&g2
    B.p=and(gp[0].p,gp[1].p);    // p1&p2=P
    B.c=or(gp[1].g,b.mem2);    // g2|(p2&c)=C
    B.mem2=and(gp[1].p,c);    // p2&c
    B.mem3='0';    // neutilizat
    B.s='0';    // neutilizat
    return(B);    //trimite b
}

int carry_lookahead (long int _a ,long int _b ,long int _c)//gest. etapelor
{
    int i,j=0,test=1;
    char a[marime+1],b[marime+1],c[marime+1],C[marime+1],s[marime+1];
    stare ch, gp1[marime],gp2[marime/2],gp3[2],gp4;

    init_stare(gp1, gp2, gp3, &gp4);
    genereaza_sir(_c, C);    // creeaza un sir pentru C nou
    genereaza_sir(_c, c);    // creeaza un sir pt. c vechi
    genereaza_sir(_a, a);    // creeaza un sir pt. a
    genereaza_sir(_b, b);    // creeaza un sir pt. b
    while(test!=0)
    {
        test=0;    // initializare test
        j++;
        // nivel 4: 1 element
        ch=element_b(gp3, gp4, c[7]);    // ultima celula B opereaza
        test=test+compara(gp4, ch);    // test de stabilitate
    }
}

```

ANEXA 2

```

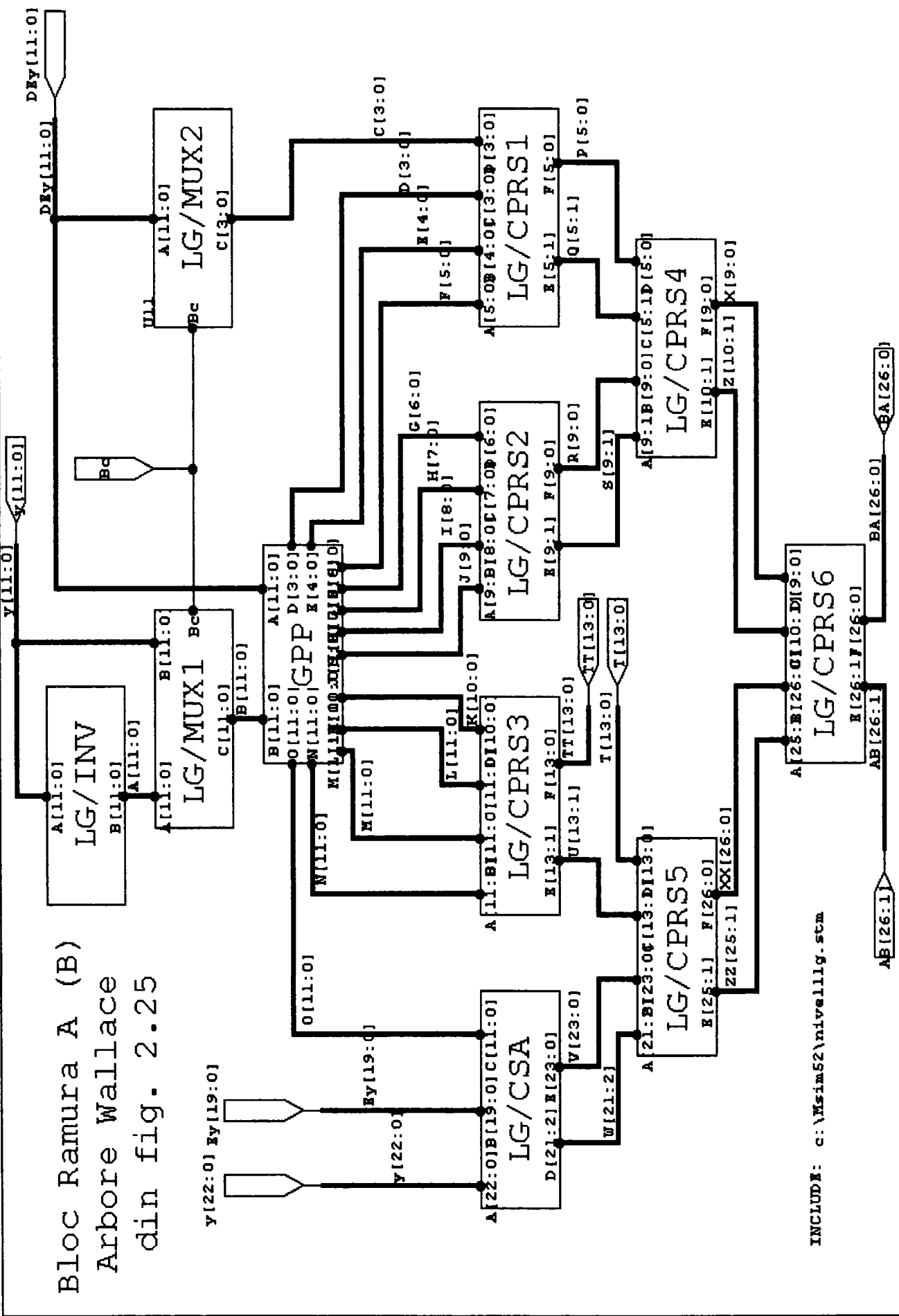
gp4=ch;                // actualizarea acestui element
C[3]=ch.c;            // memoreaza transportul
for (i=0;i<2;i++)    // nivel 3: 2 elemente
{
    ch=element_b(gp2+2*i,gp3[i],c[3+4*i]); // cele 2 celule B opereaza
    test=test+compara(gp3[i],ch); // test de stabilitate
    gp3[i]=ch;        // actualizarea acestor elemente
    C[1+4*i]=ch.c;    // memoreaza transporturile
}
for (i=0;i<4;i++)    //al doilea nivel: 4 elemente
{
    ch=element_b(gp1+2*i,gp2[i],c[1+2*i]); //cele 4 celule B opereaza
    test=test+compara(gp2[i],ch); // test de stabilitate
    gp2[i]=ch;        // actualizarea acestor elemente
    C[2*i]=ch.c;     // memoreaza transporturile
}
for (i=0;i<marime;i++) // nivel 1: 8 elemente
{
    ch=element_a(a[i],b[i],c[i],gp1[i]); //cele 8 celule A opereaza
    test=test+compara(gp1[i],ch); // test de stabilitate
    gp1[i]=ch;        // actualizarea acestor elemente
    if (c[i]!=C[i])test=1; // test de stabilitate a transp.
    c[i]=C[i];        // actualizarea transp.
}
}
return(j);            // trimite numarul de etape
}

void main()
{
double max;
unsigned long int a,b,c,i;
int t,t_max=0;
max=rang(2,marime); // calcul numar maxim in functie de numar de biti
for(a=0;a<max;a++) // bucle pt. toate combinatiile a,b,c
{
    for(b=0;b<max;b++)
    {
        for(c=0;c<2;c++)
        {
            t=carry_lookahead(a,b,c); // calculeaza timpul de calcul
            if(t_max<t) // daca s-a depasit valoarea maxima
            {
                i=0; // initializeaza i
                t_max=t; // memoreaza aceasta valoare
                printf("carry_lookahead:\n"); // tipareste
                printf("timp de propagare:%d valori:\n",t_max);
            }
            if(t_max==t) // pt. acest timp de propagare
            {
                printf("%ld,%ld,%ld\t",a,b,c); // tipareste valorile a, b si c
                i++; // incrementeaza i
            }
        }
    }
}
printf("\ntimpul maxim este %d timpi de propagare pentru %ld valori ale lui a,b,c",t_max,i);
getch(); // tipareste rezultatul final si asteapta
}

```

ANEXA 3

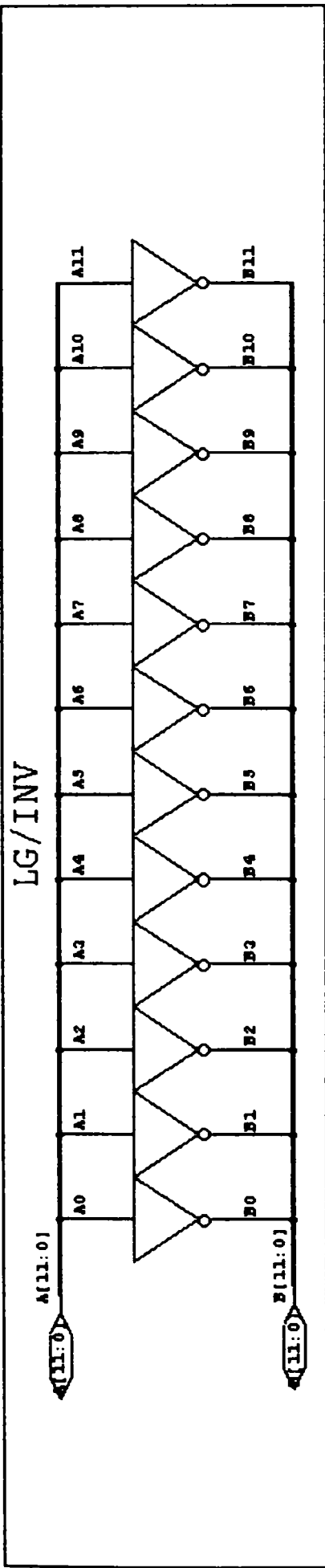
Bloc Ramura A (B)
Arbore Wallace
din fig. 2.25



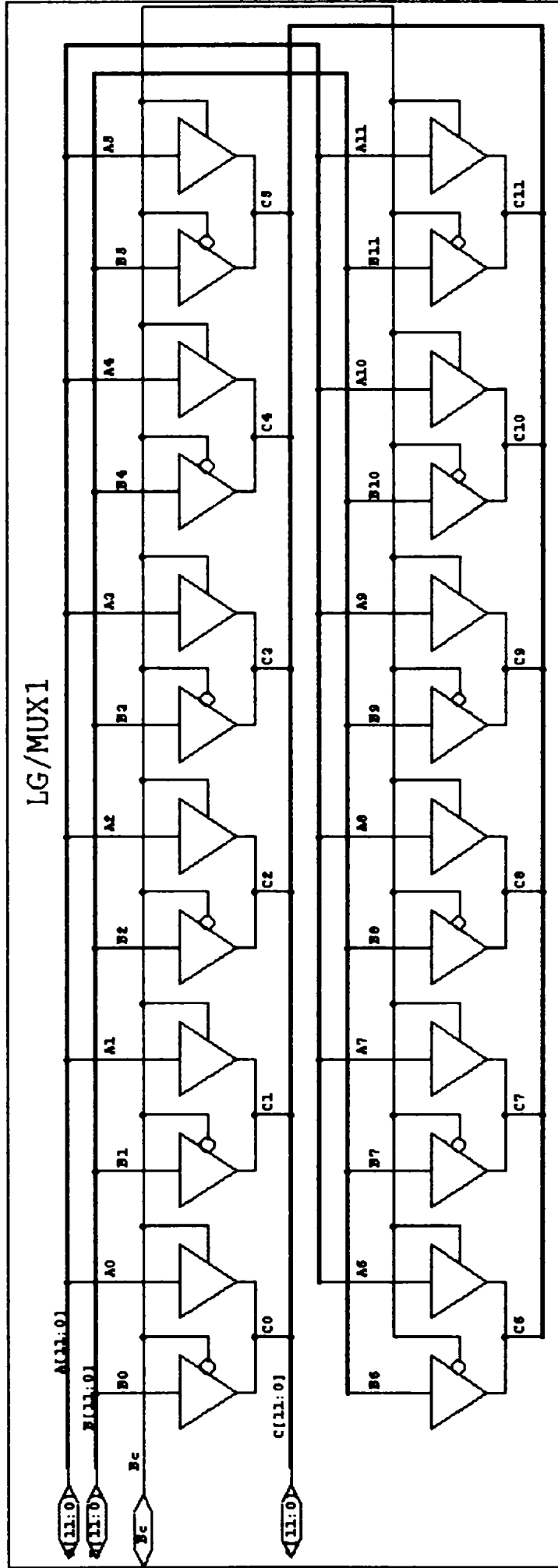
INCLUDE: c:\msim52\nivell1g.stm

ANEXA 3

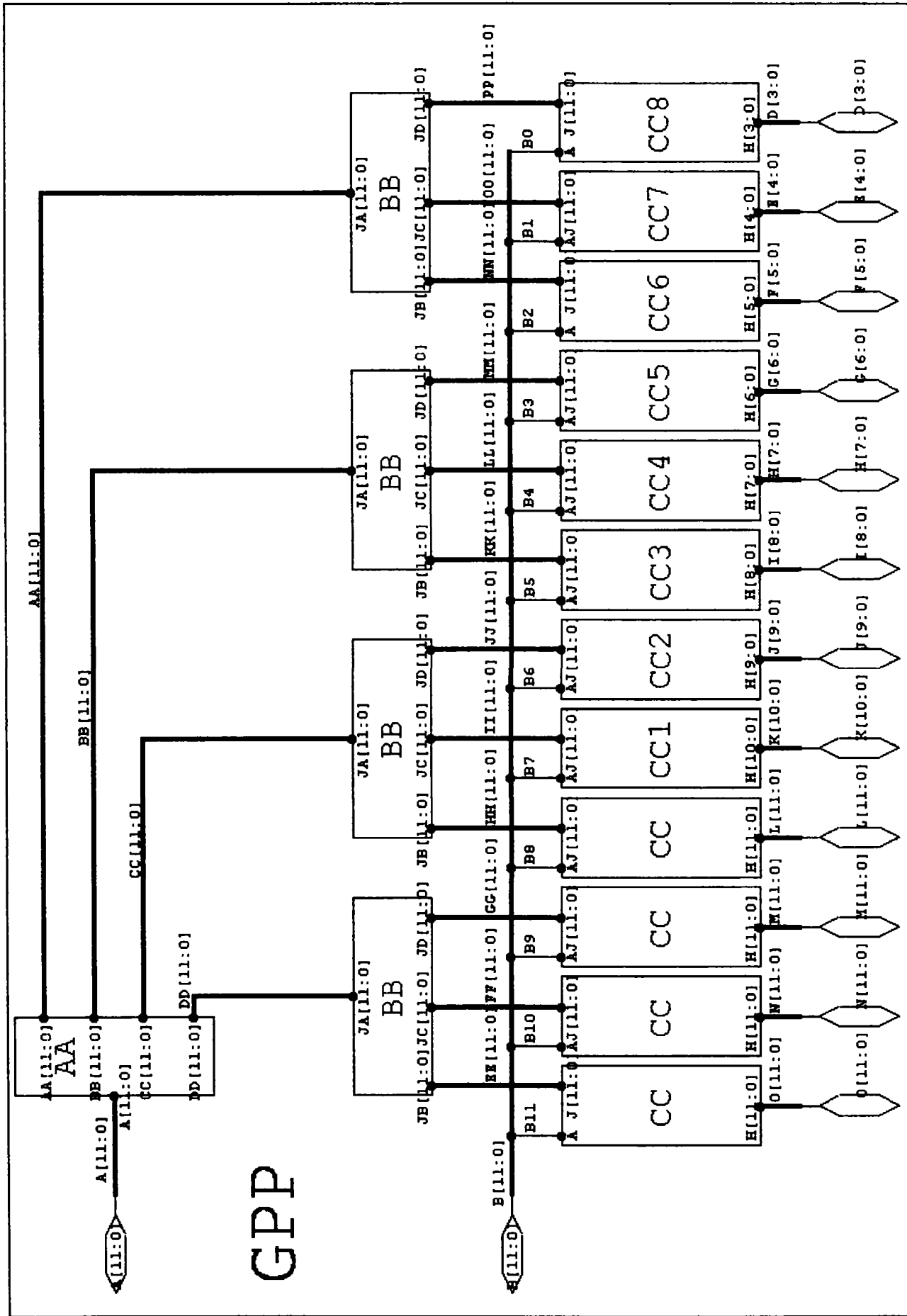
LG/INV



LG/MUX1

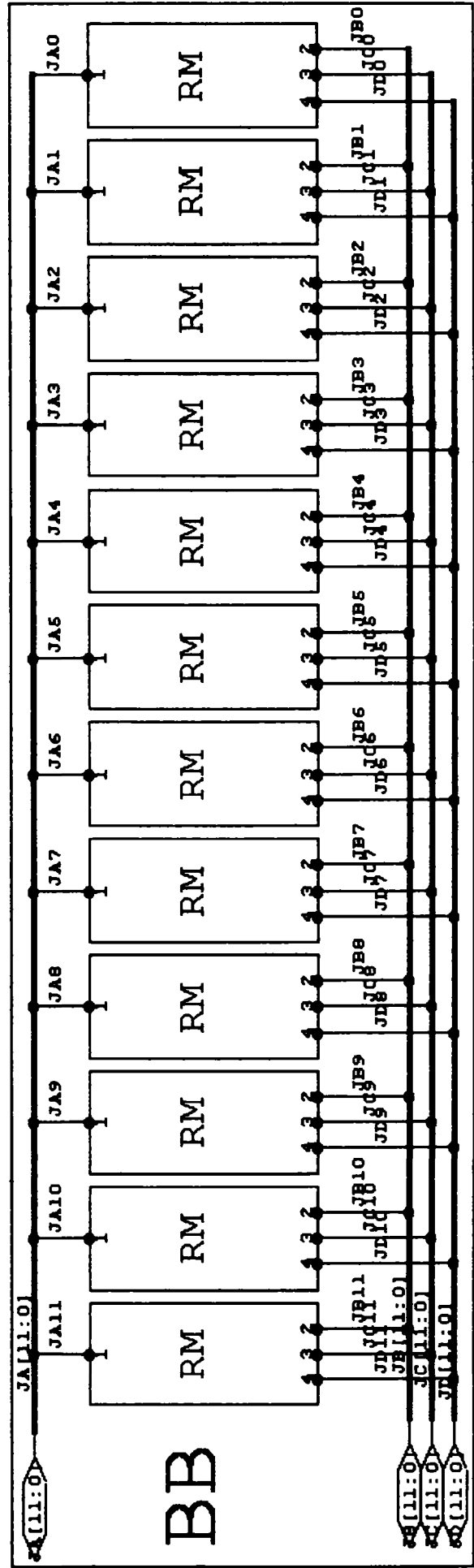
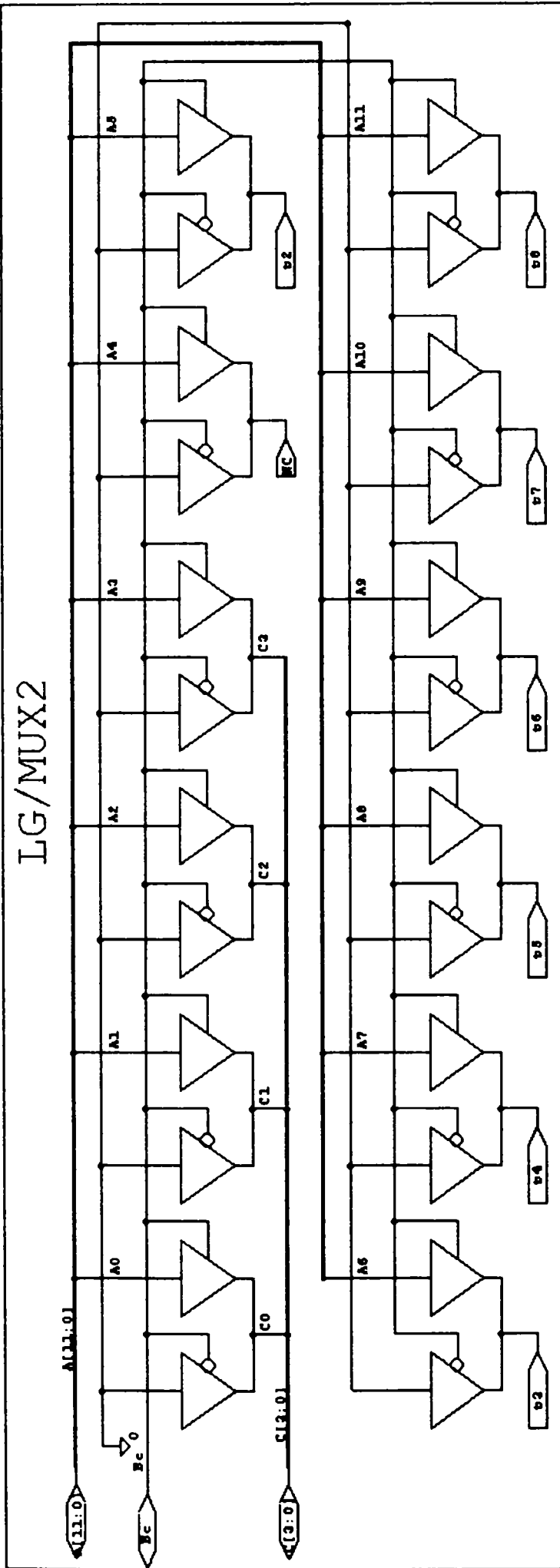


ANEXA 3

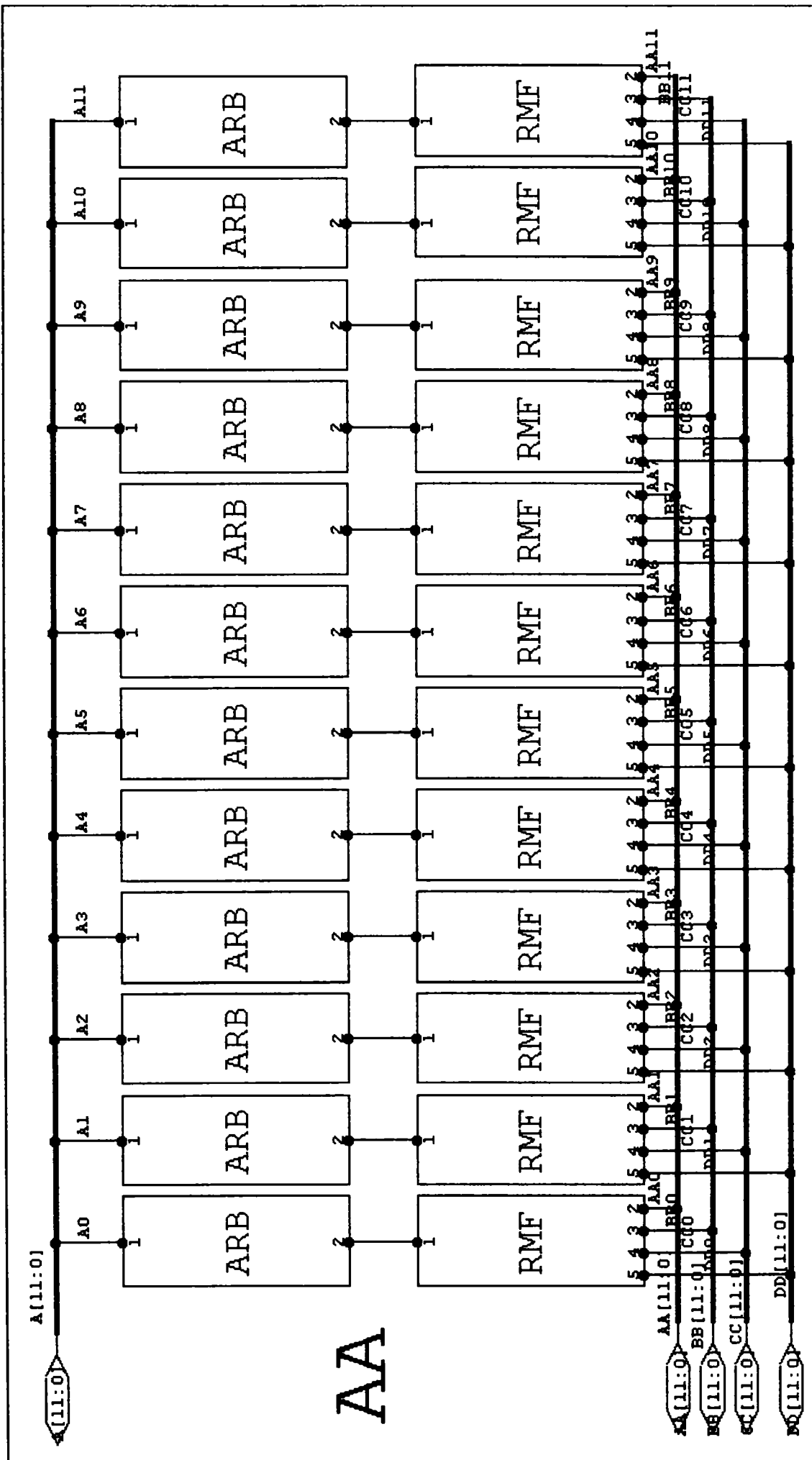


ANEXA 3

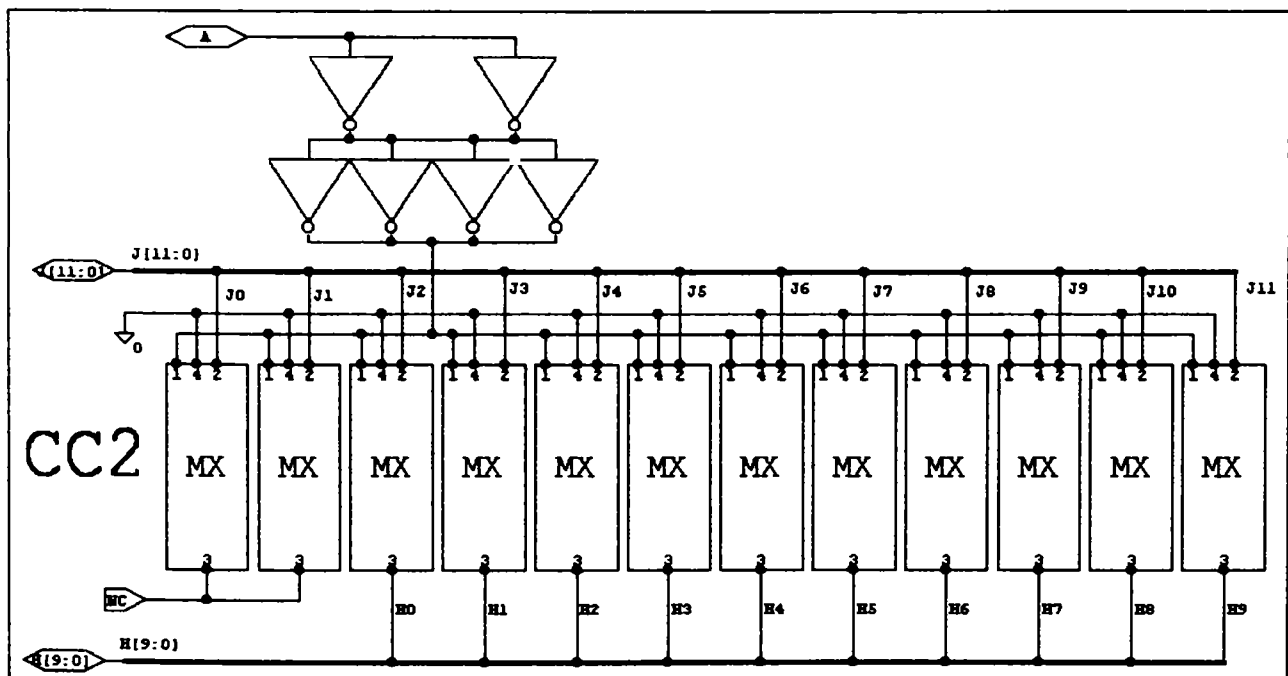
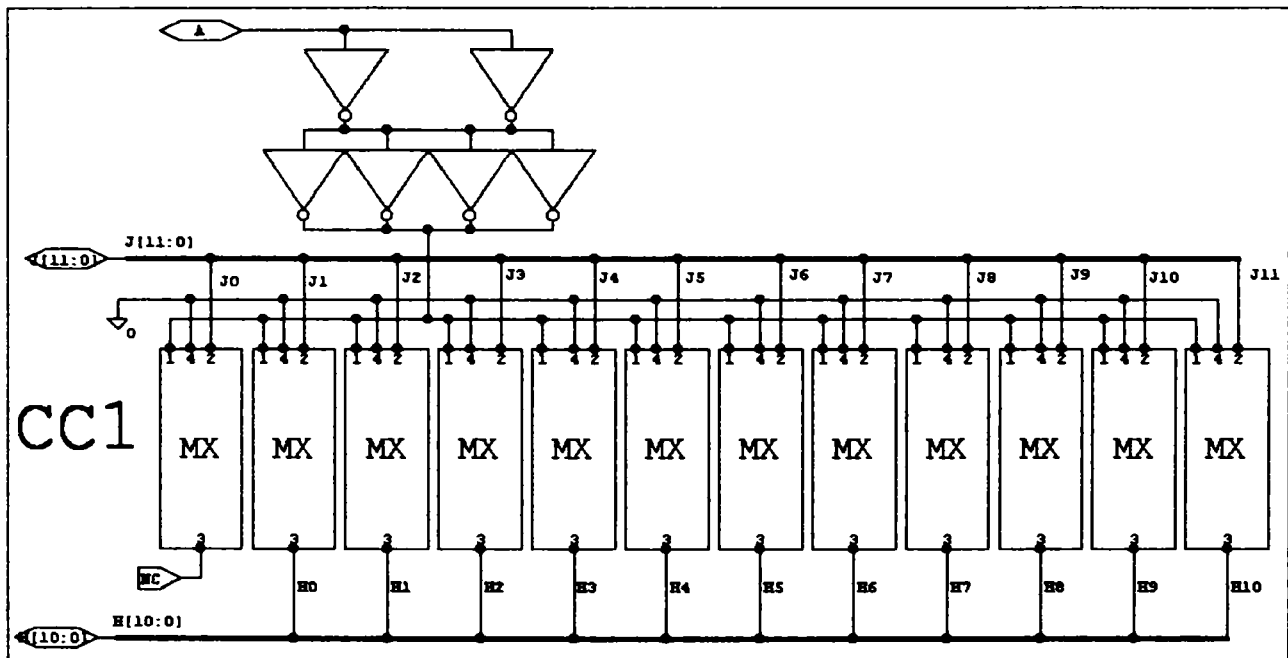
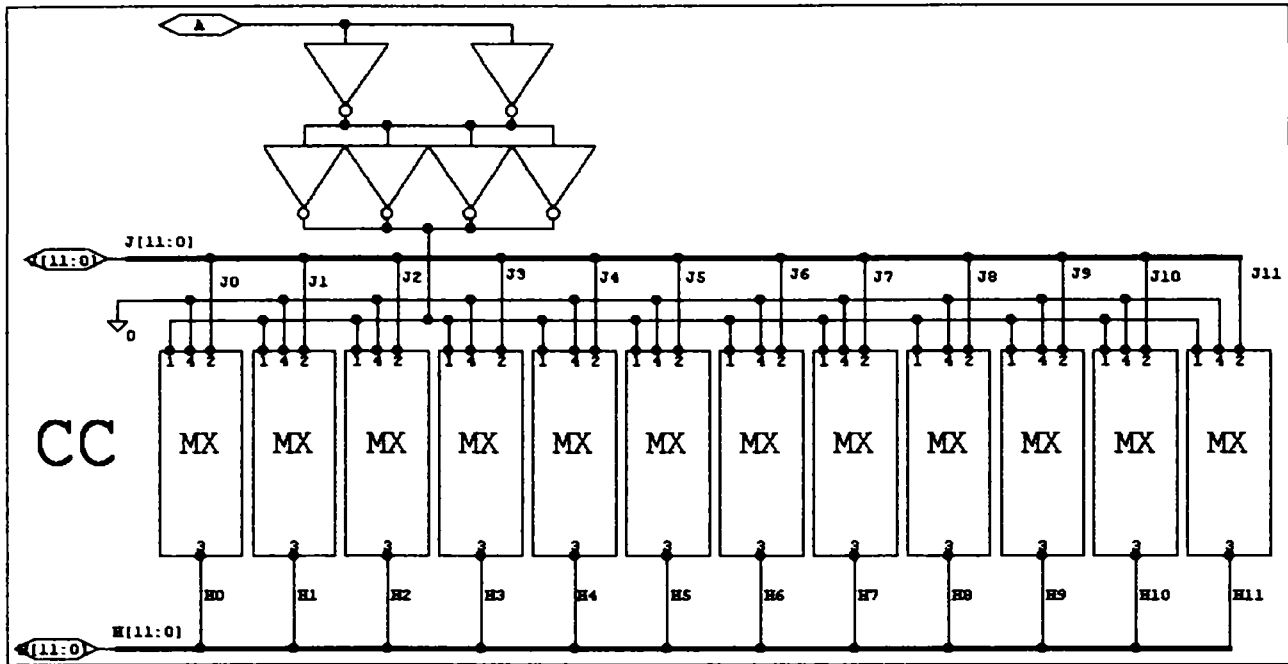
LG/MUX2

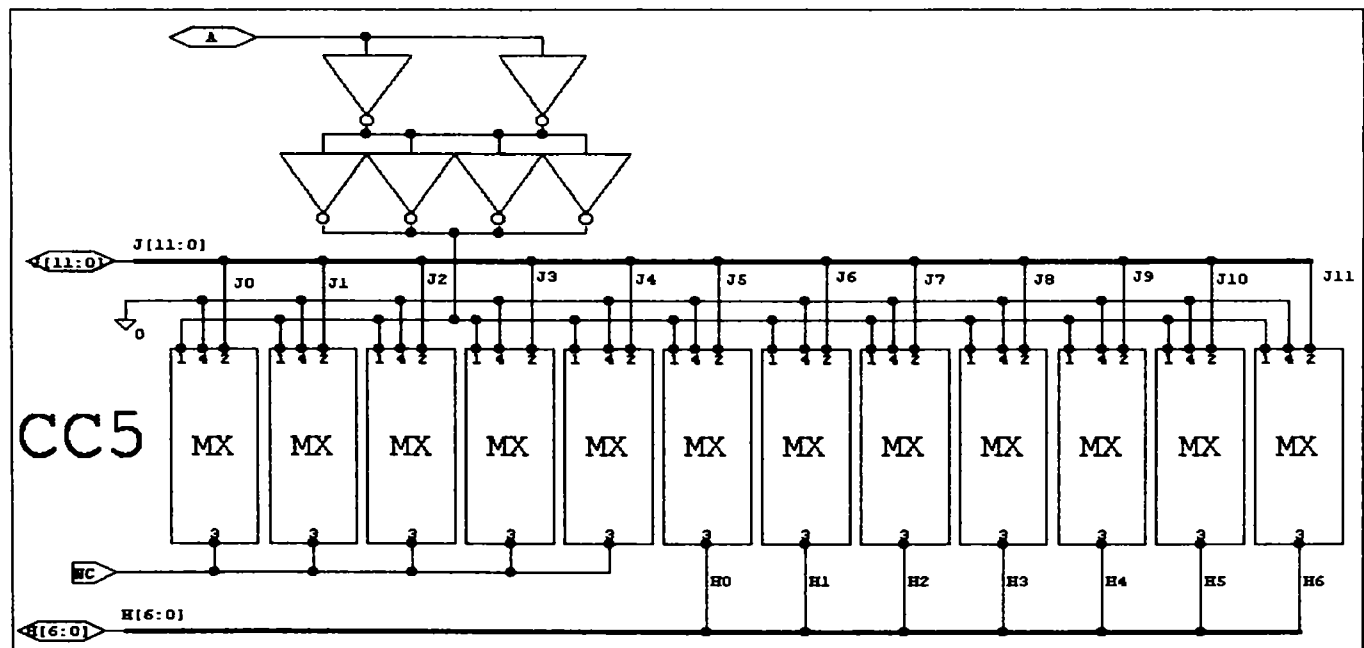
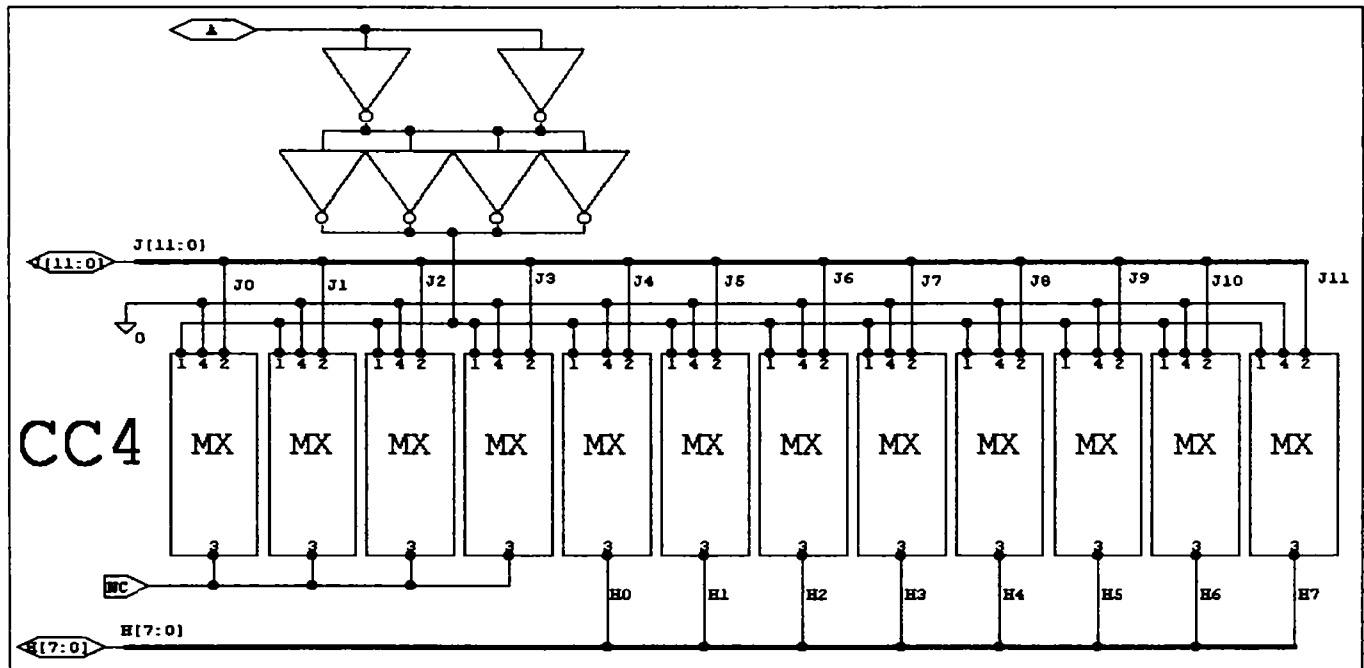
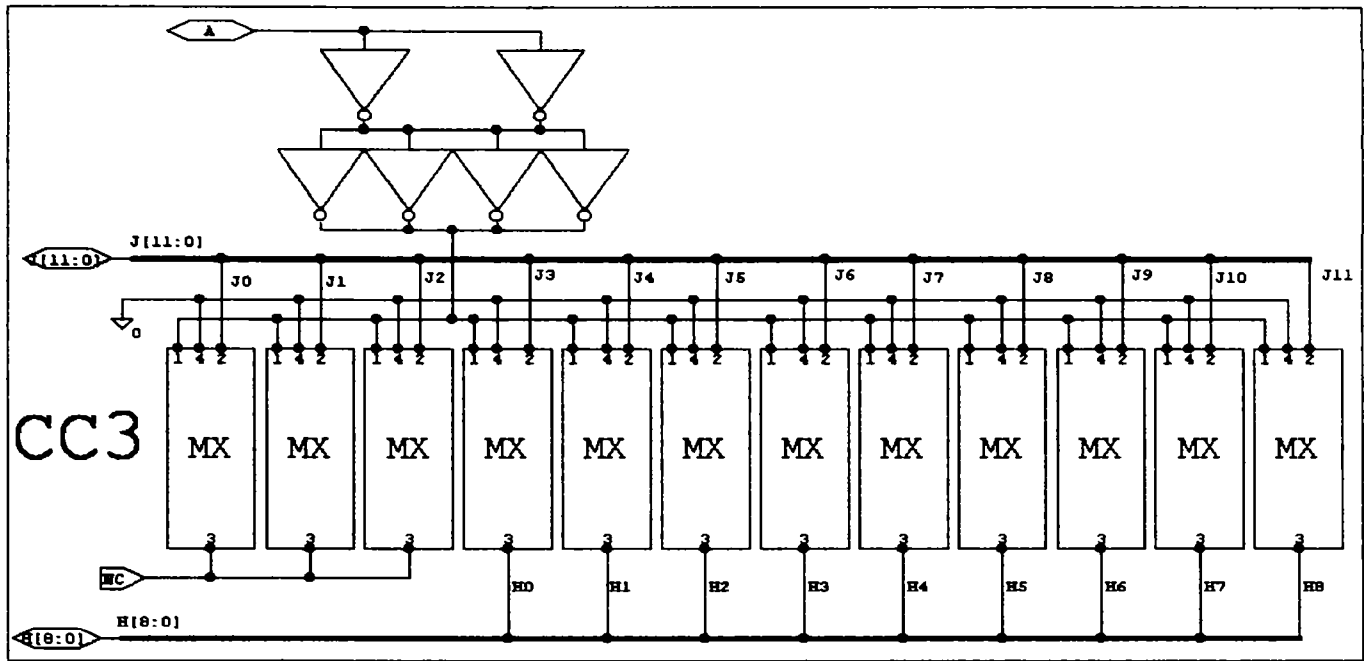


ANEXA 3

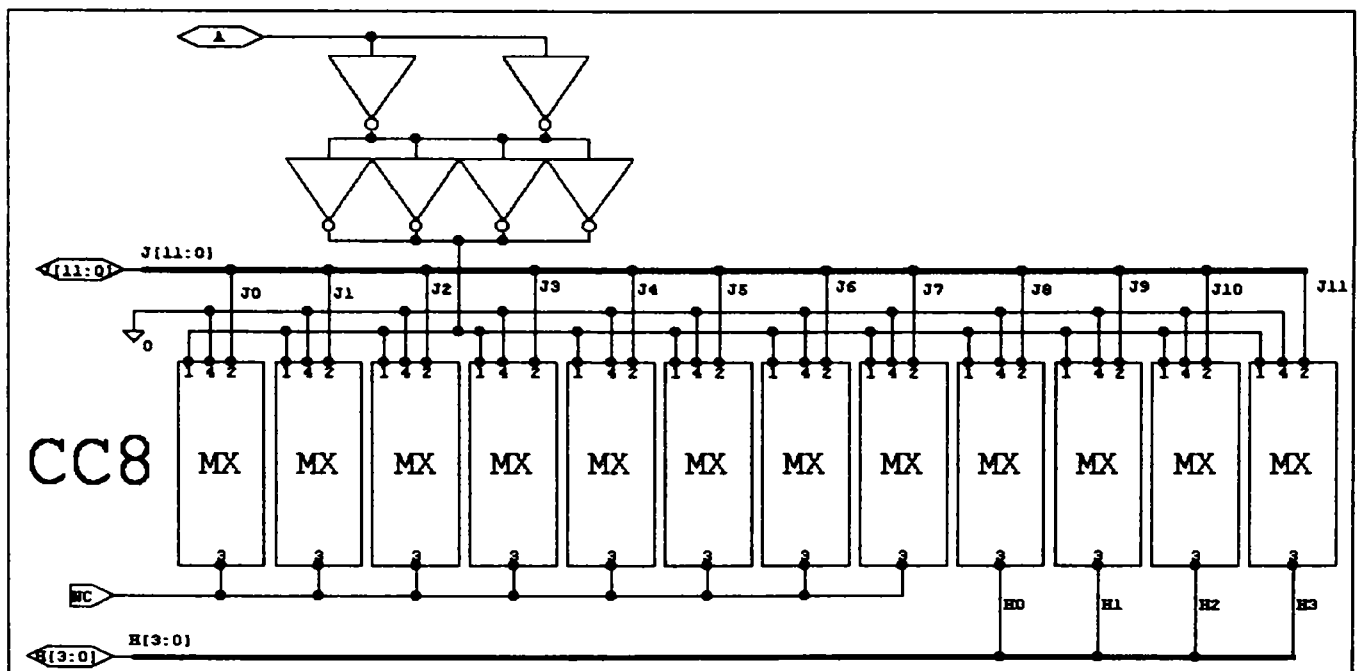
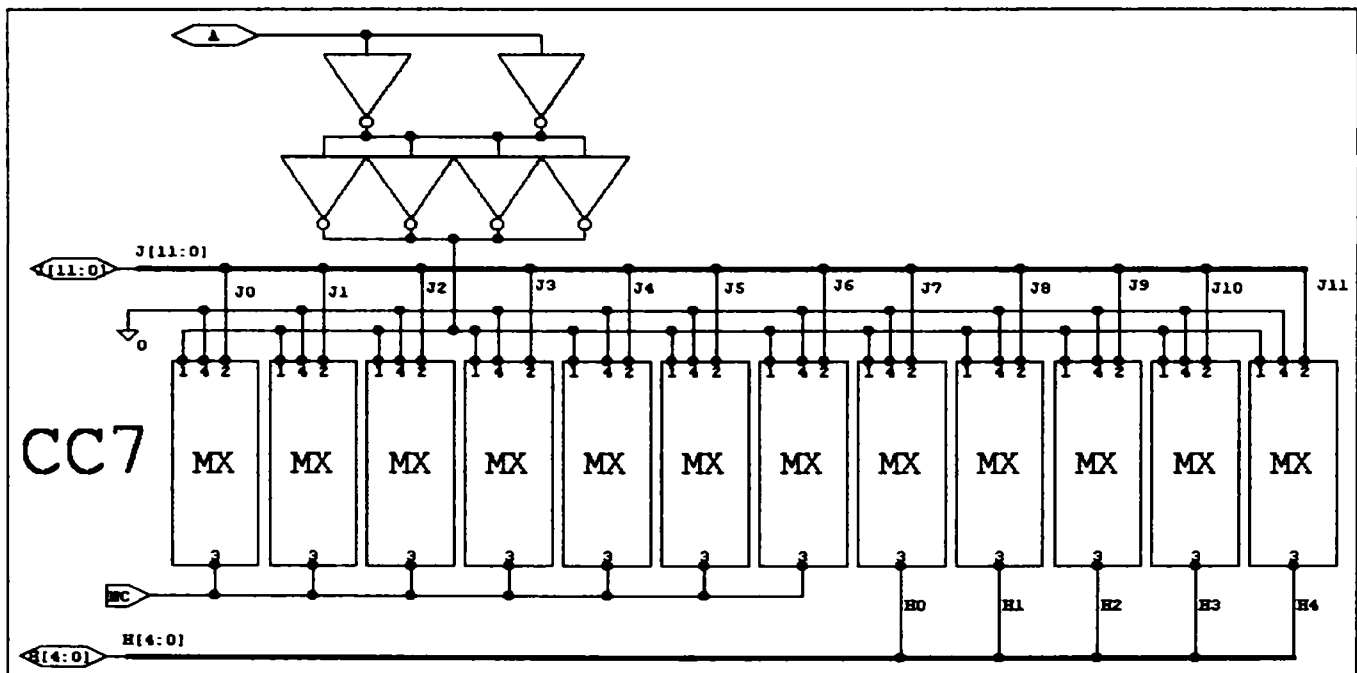
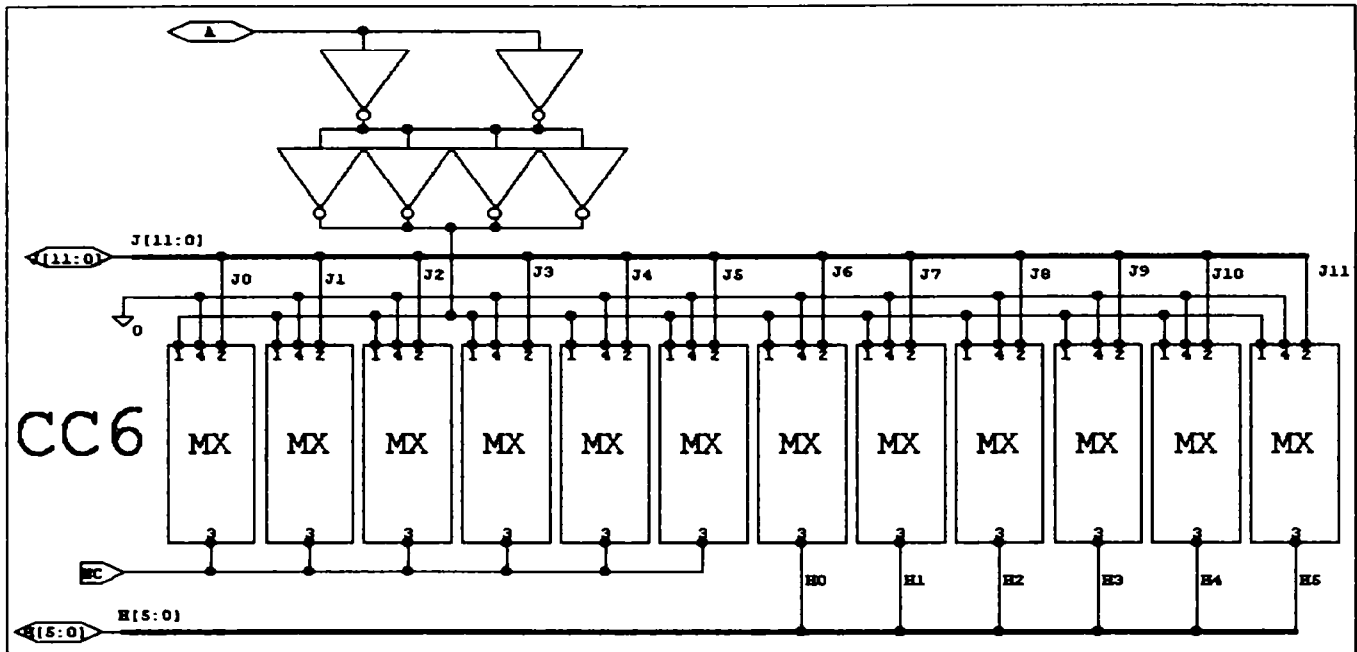


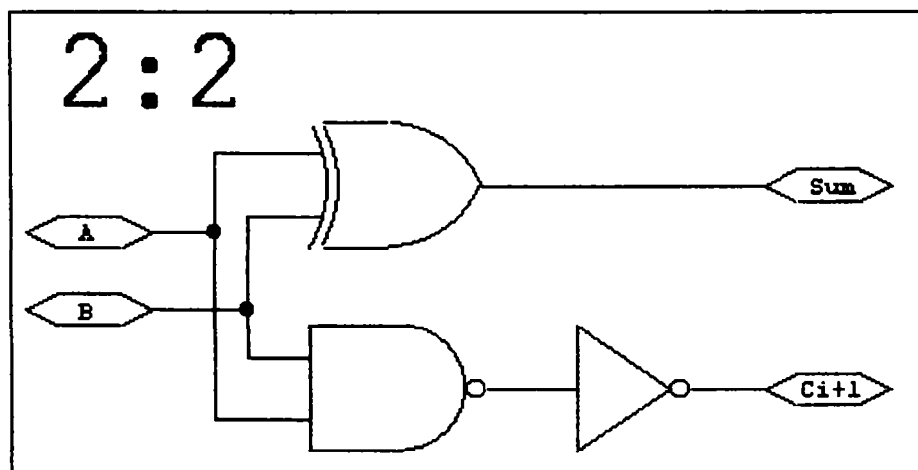
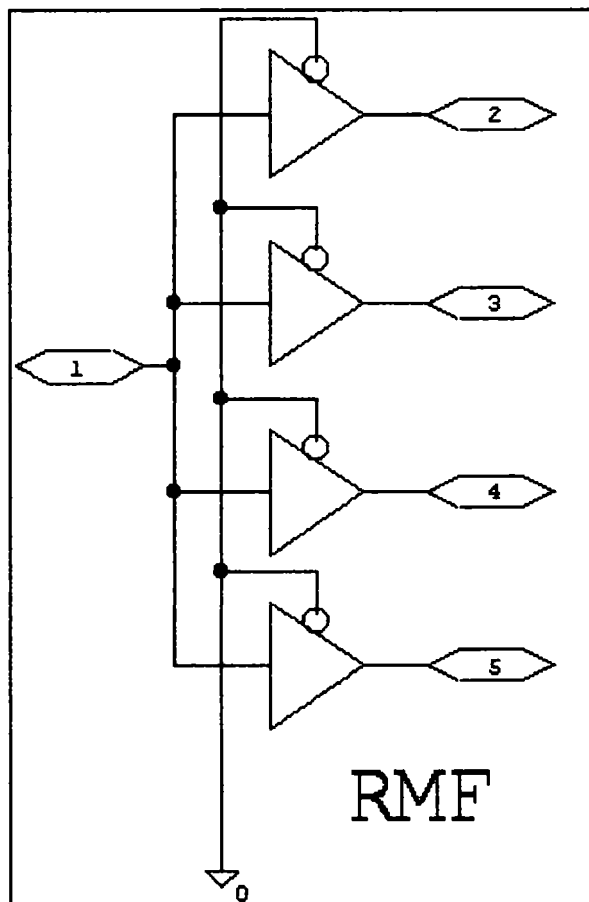
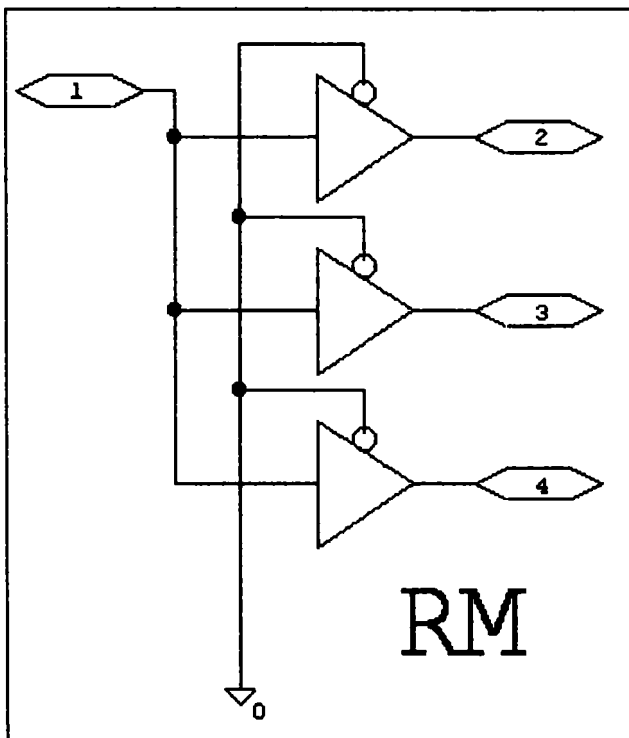
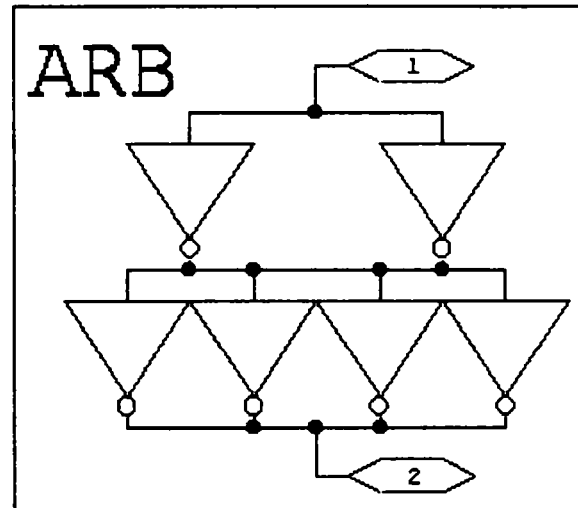
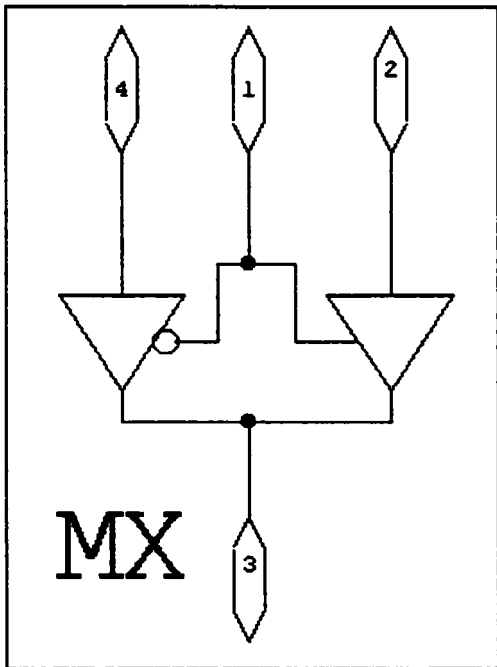
ANEXA 3



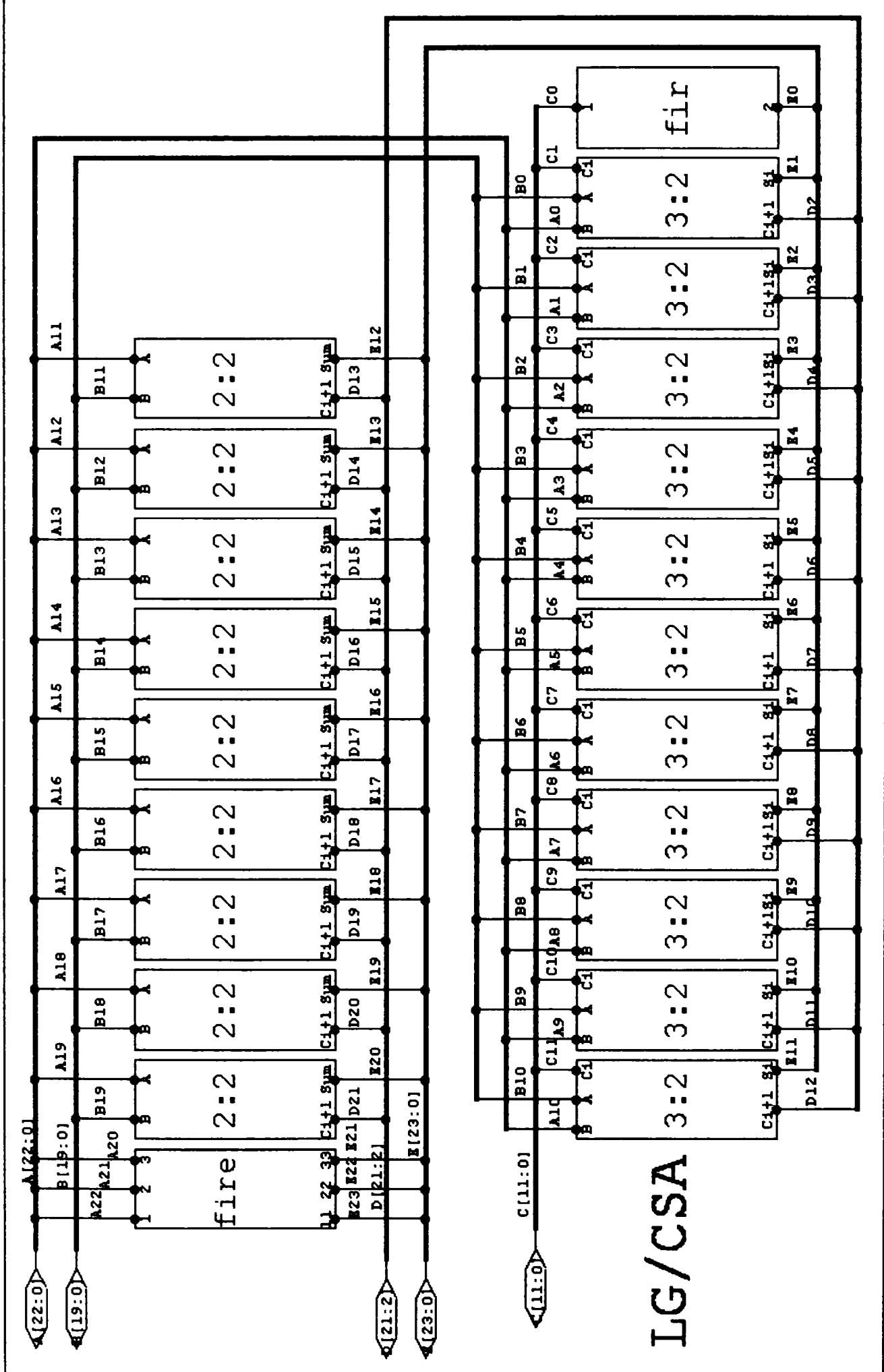


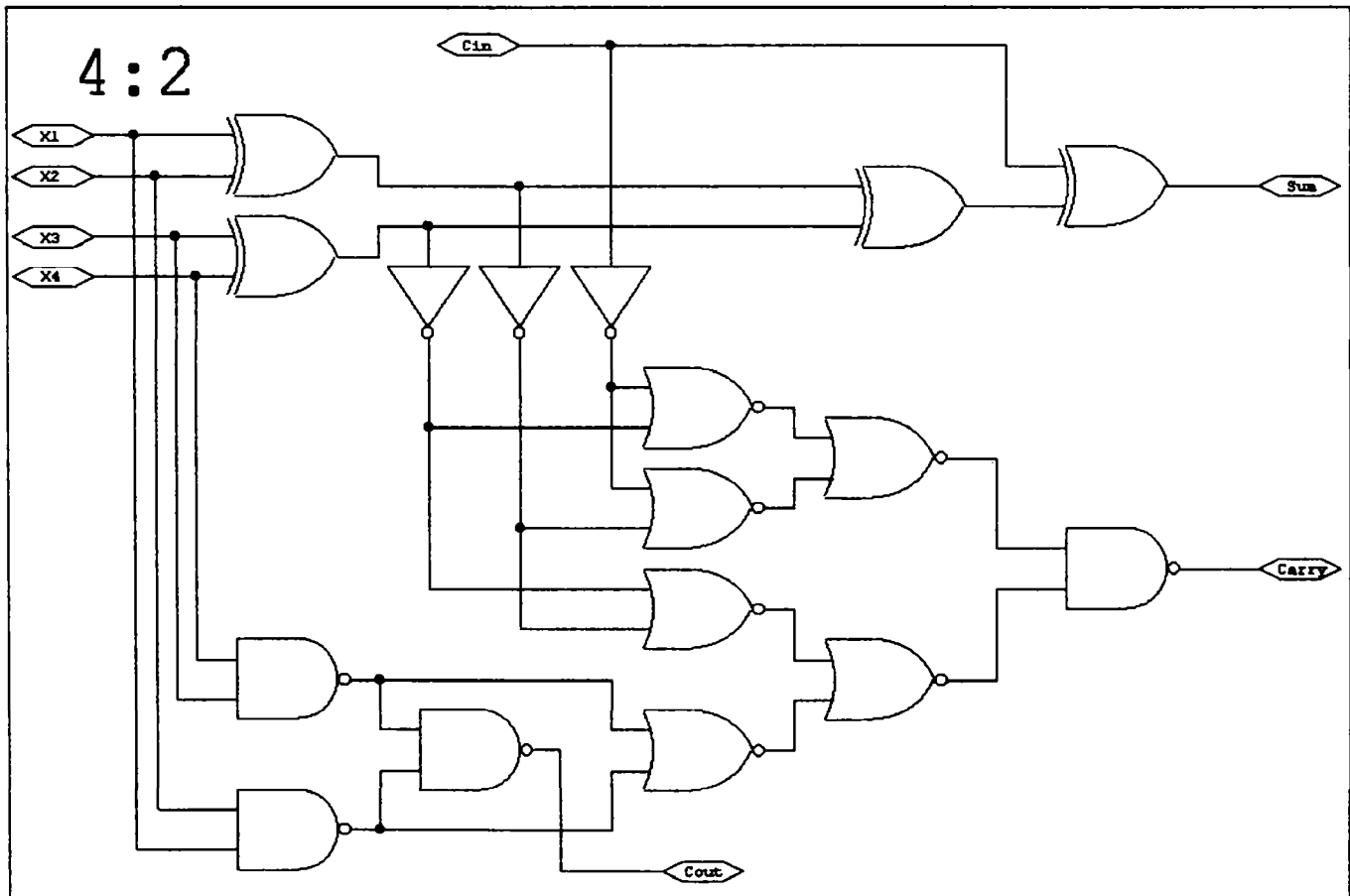
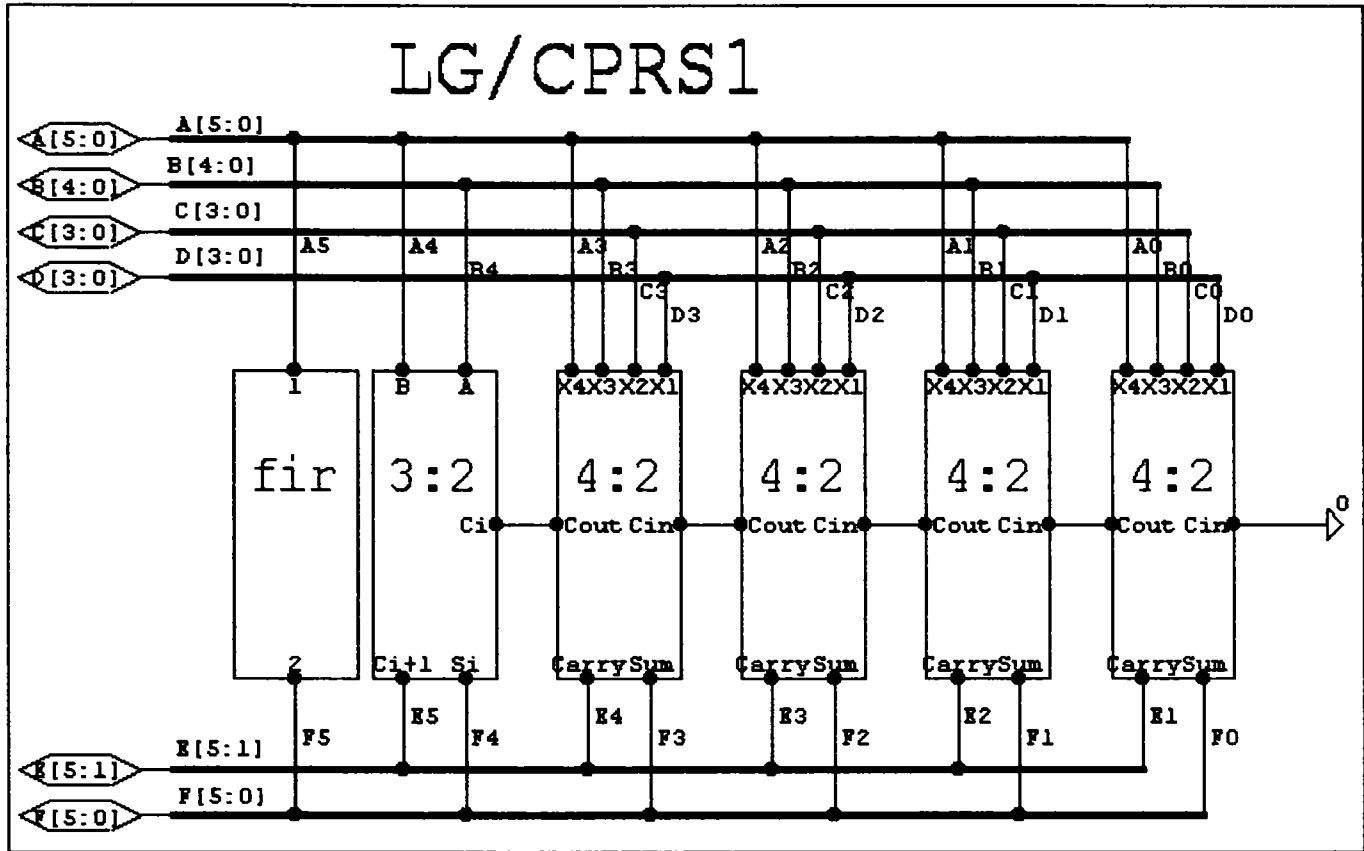
ANEXA 3



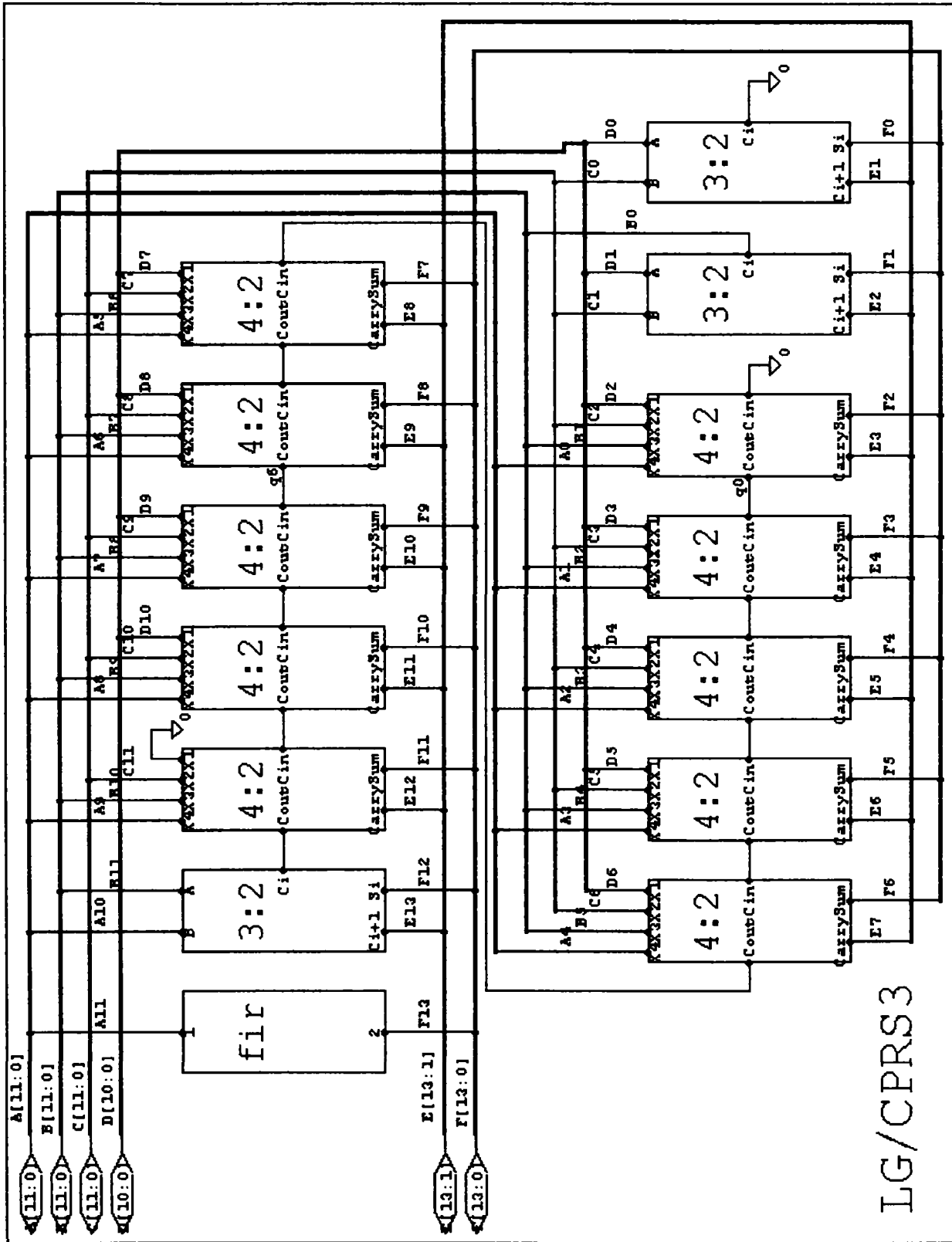


ANEXA 3

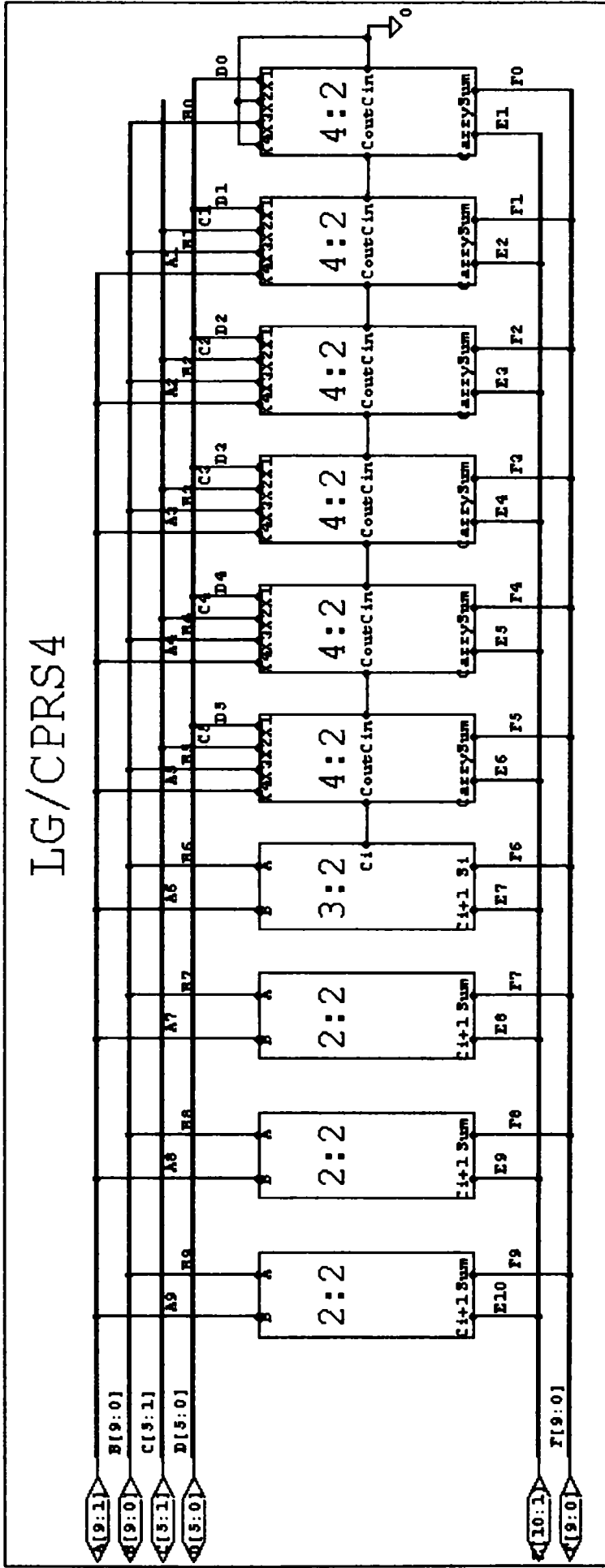




ANEXA 3

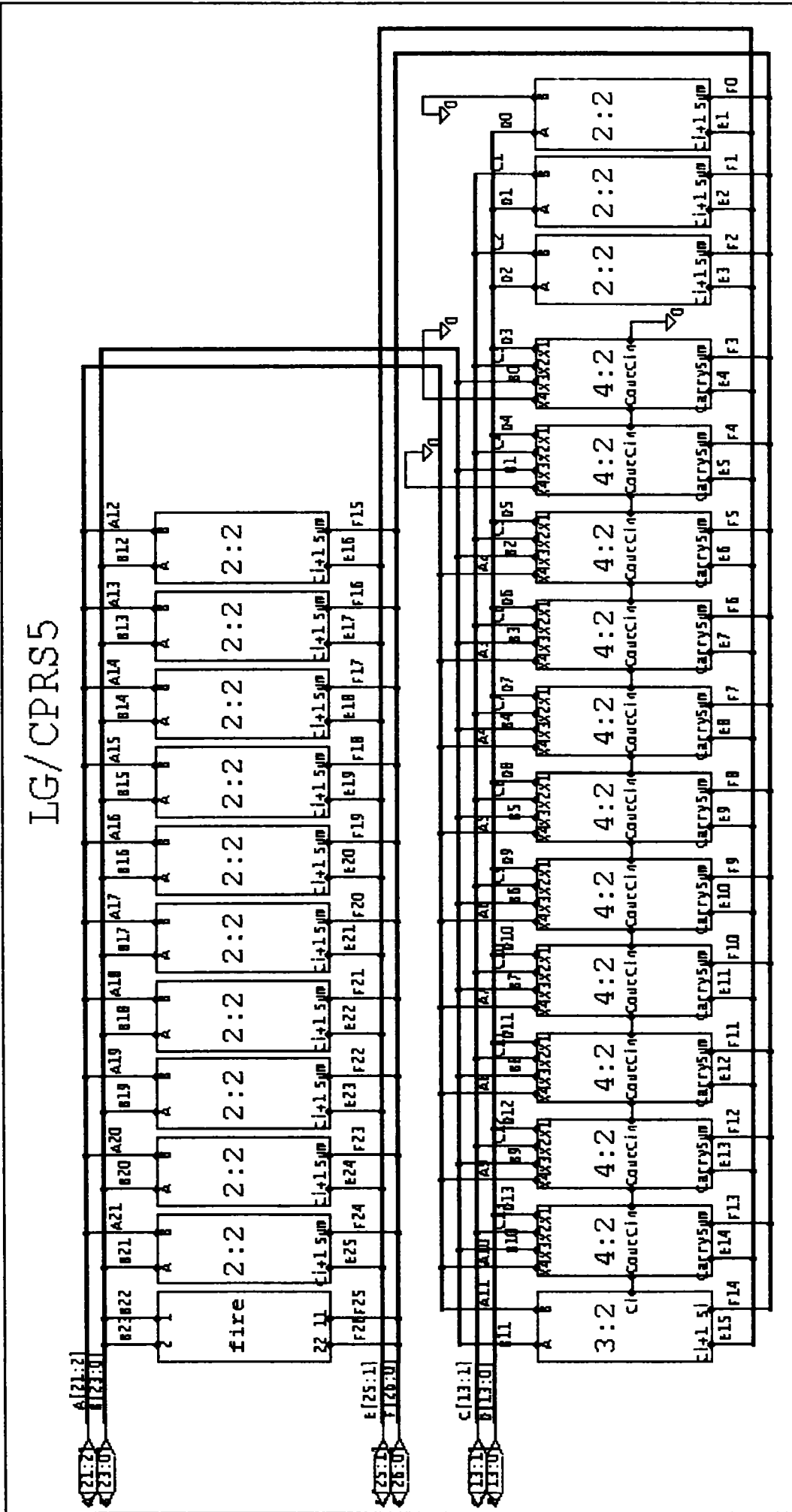


ANEXA 3

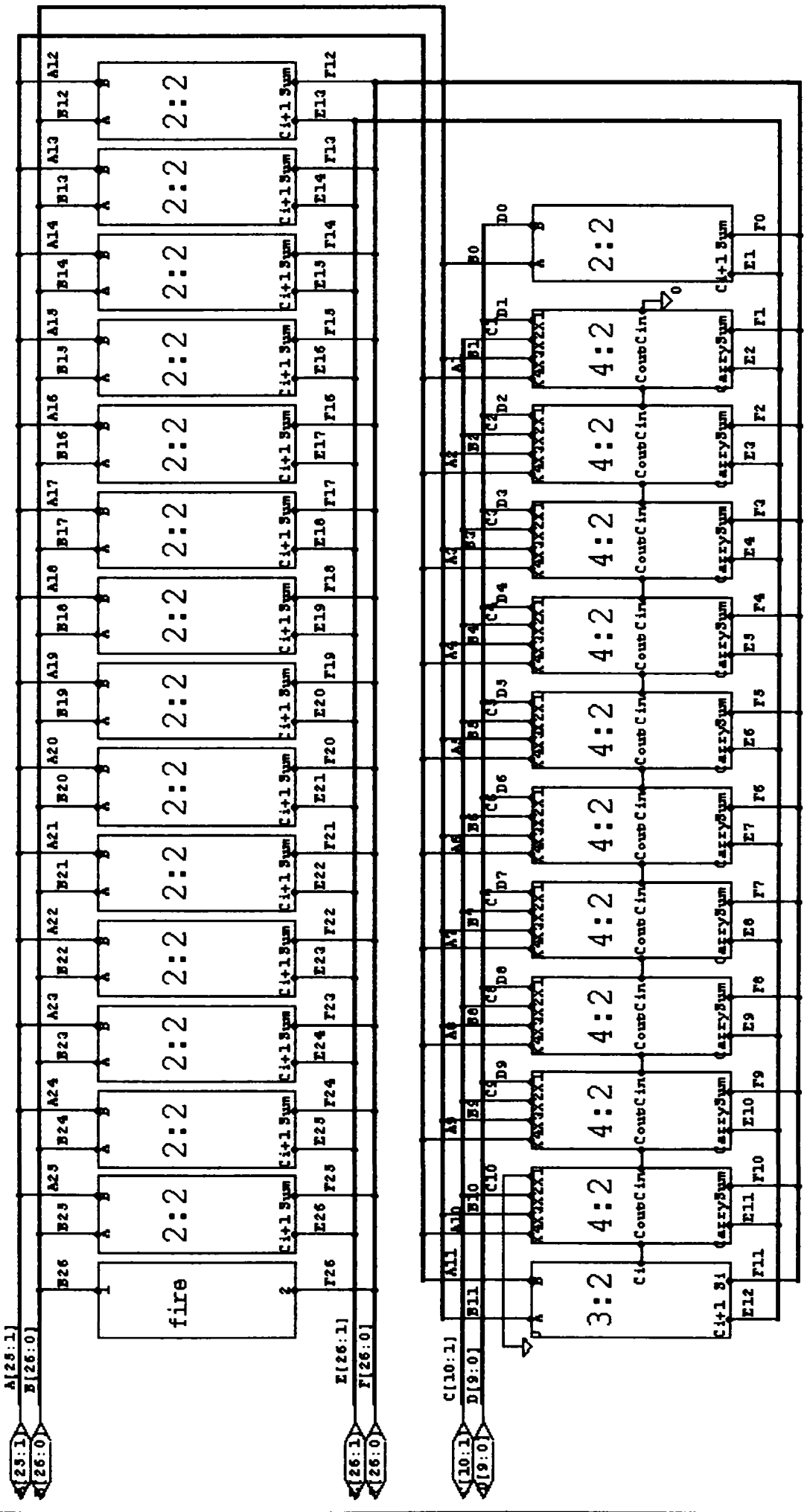


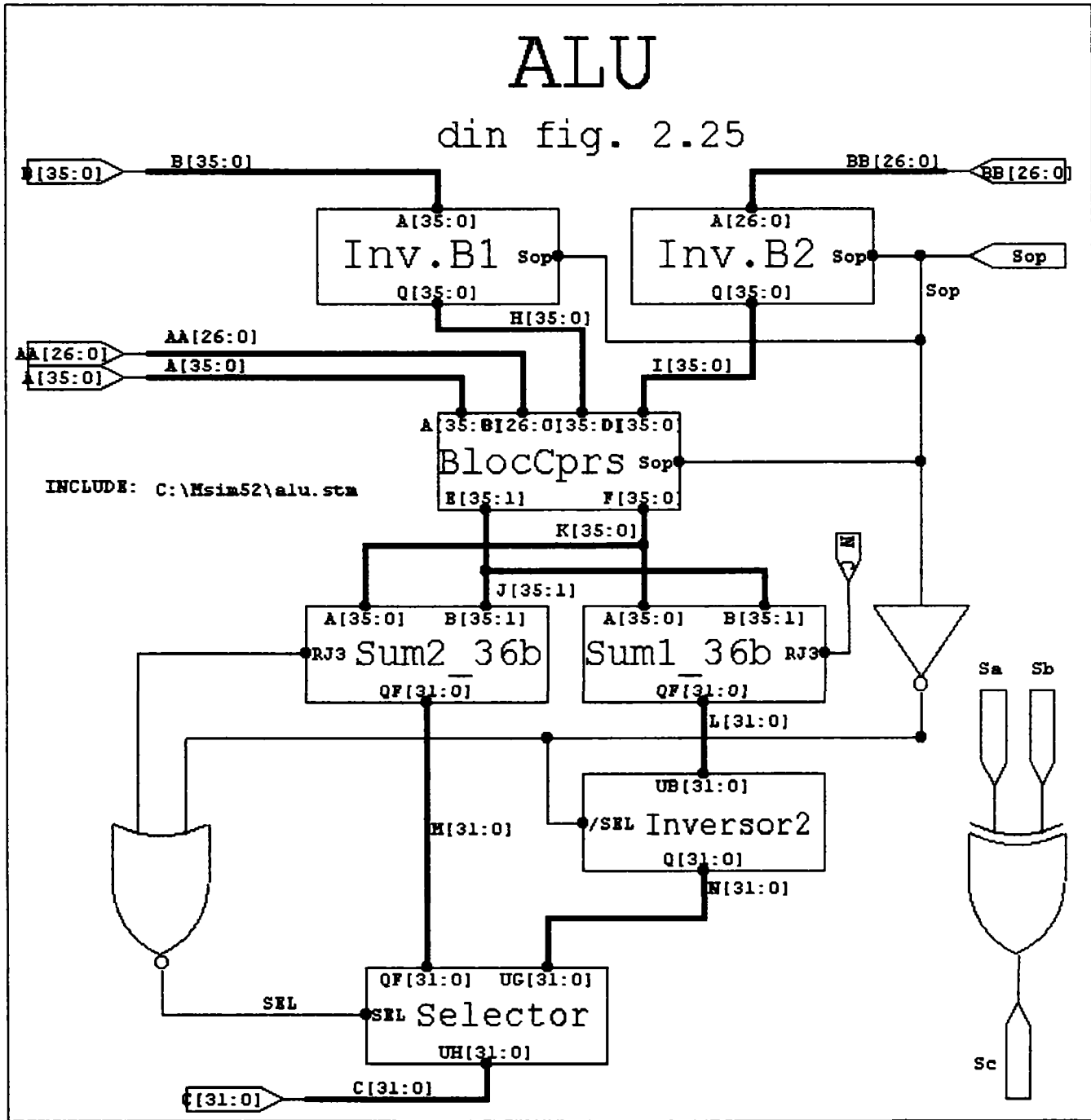
ANEXA 3

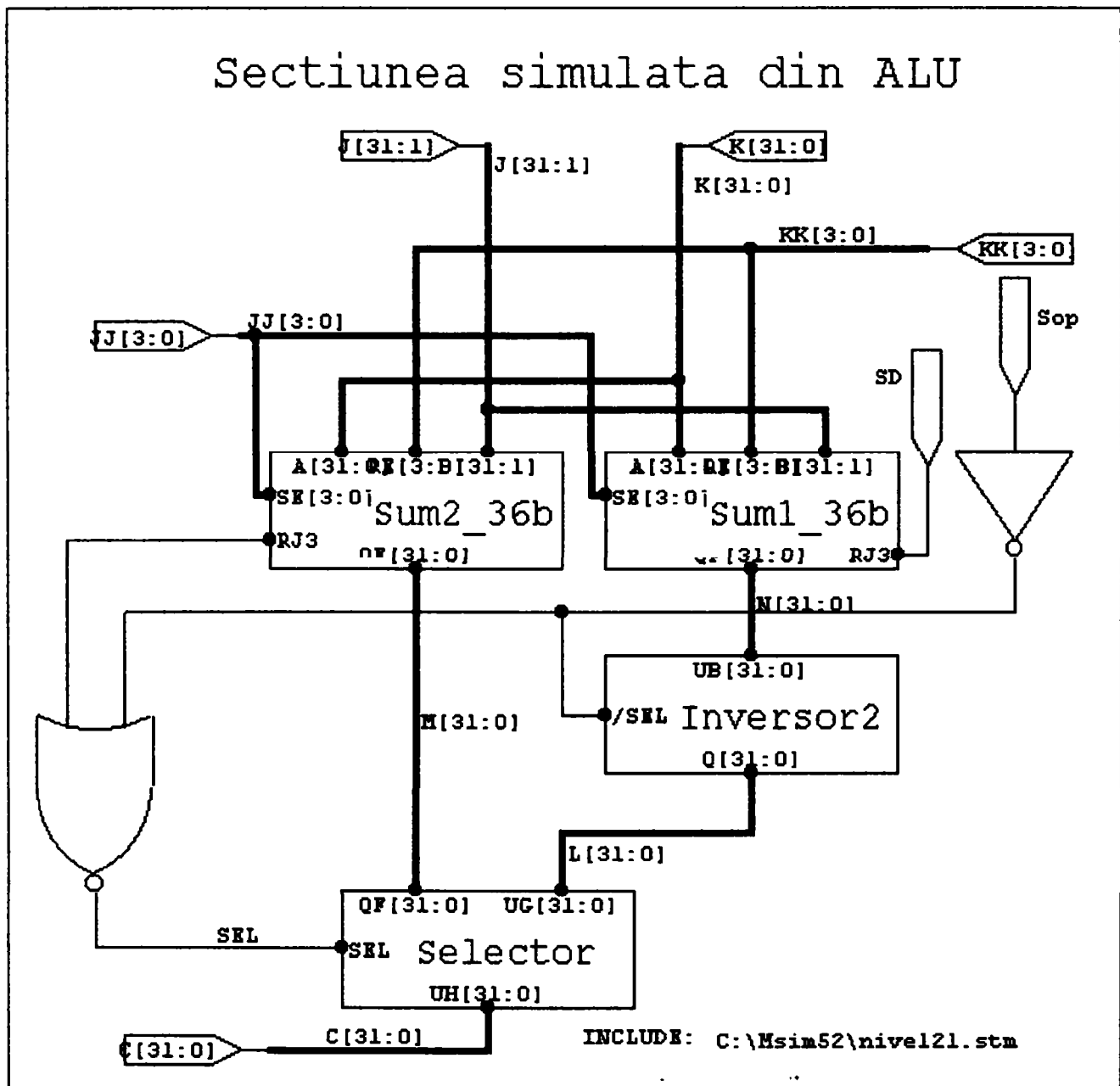
IG/ CPRS5



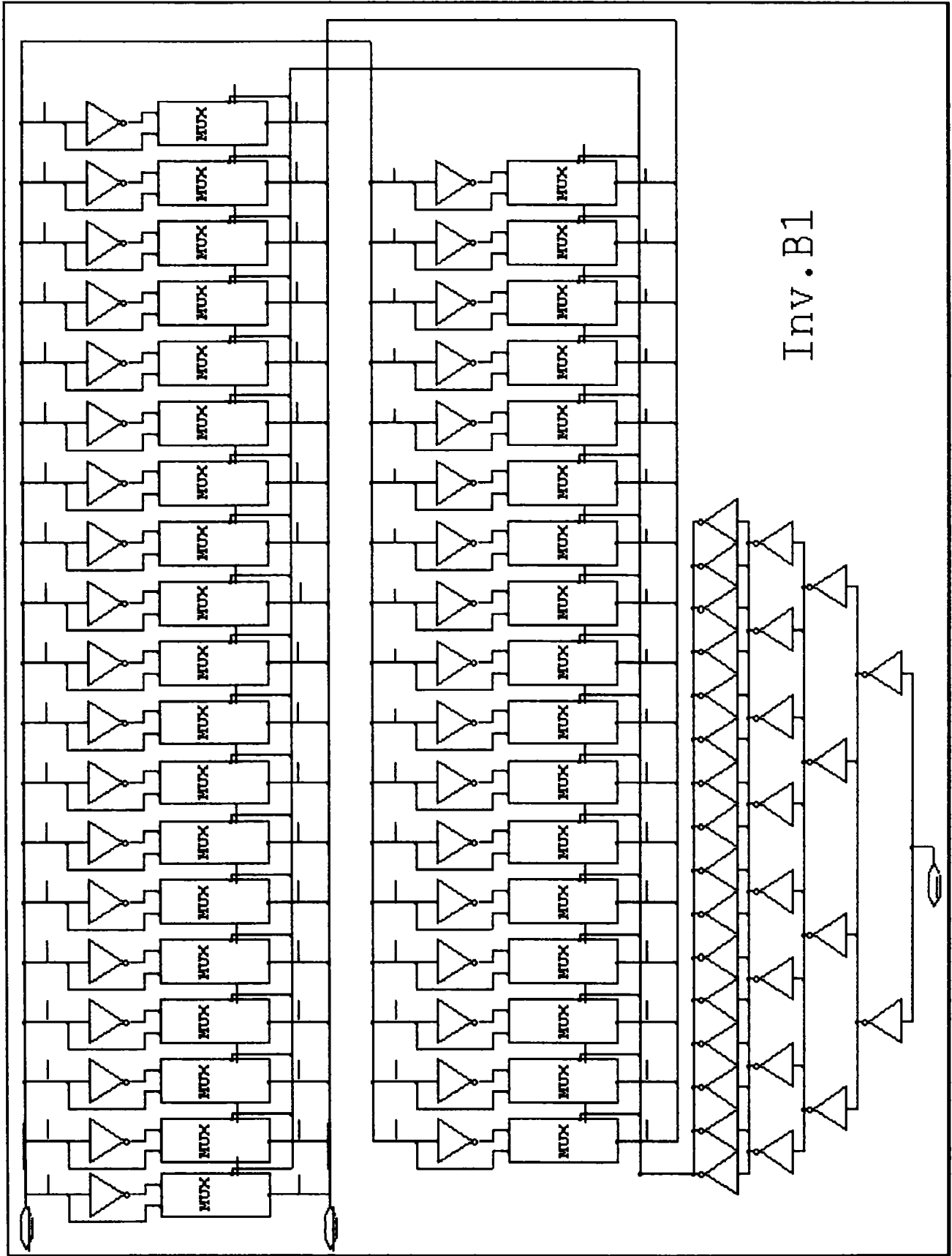
IG/CPRS6



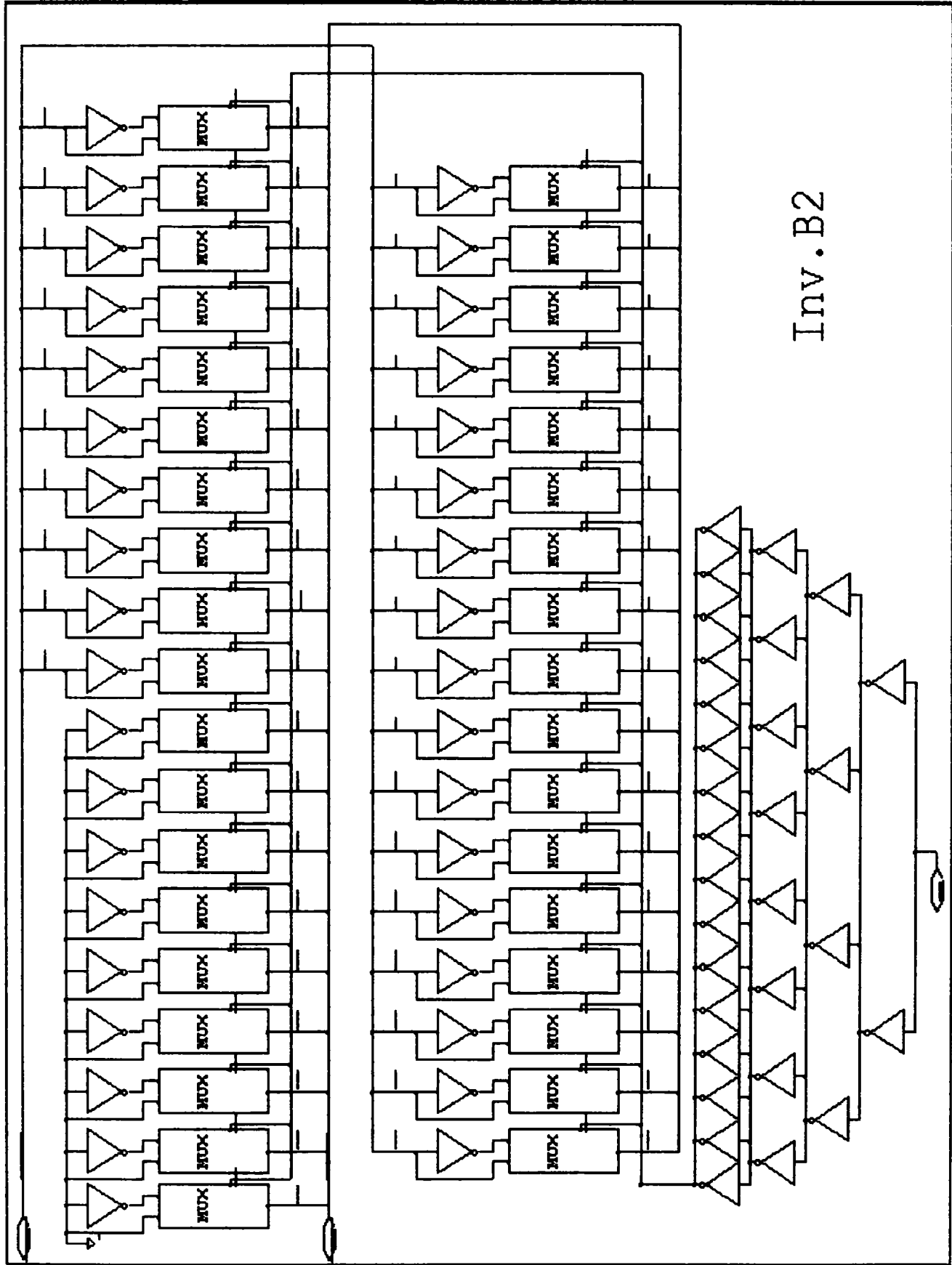




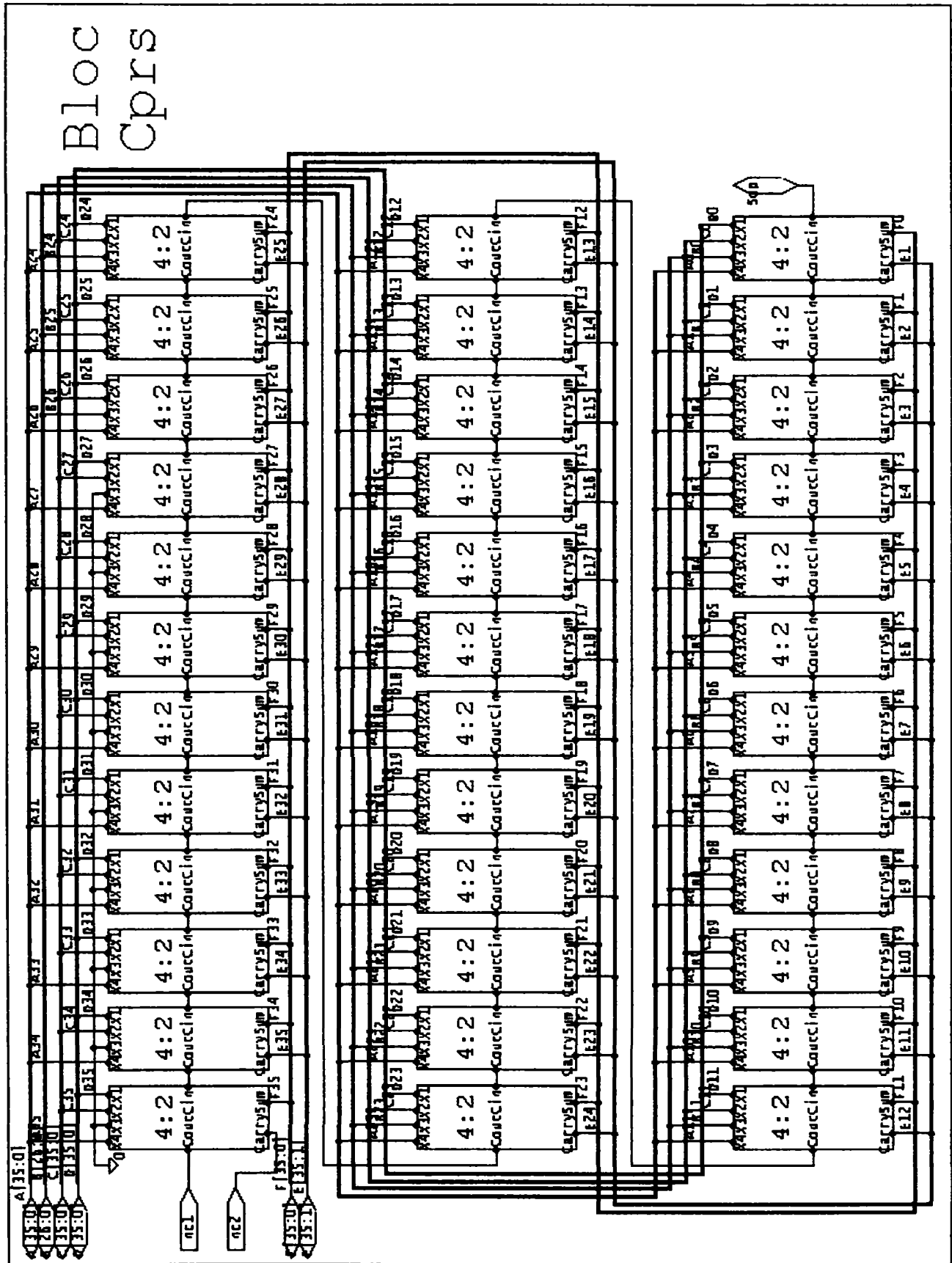
ANEXA 3

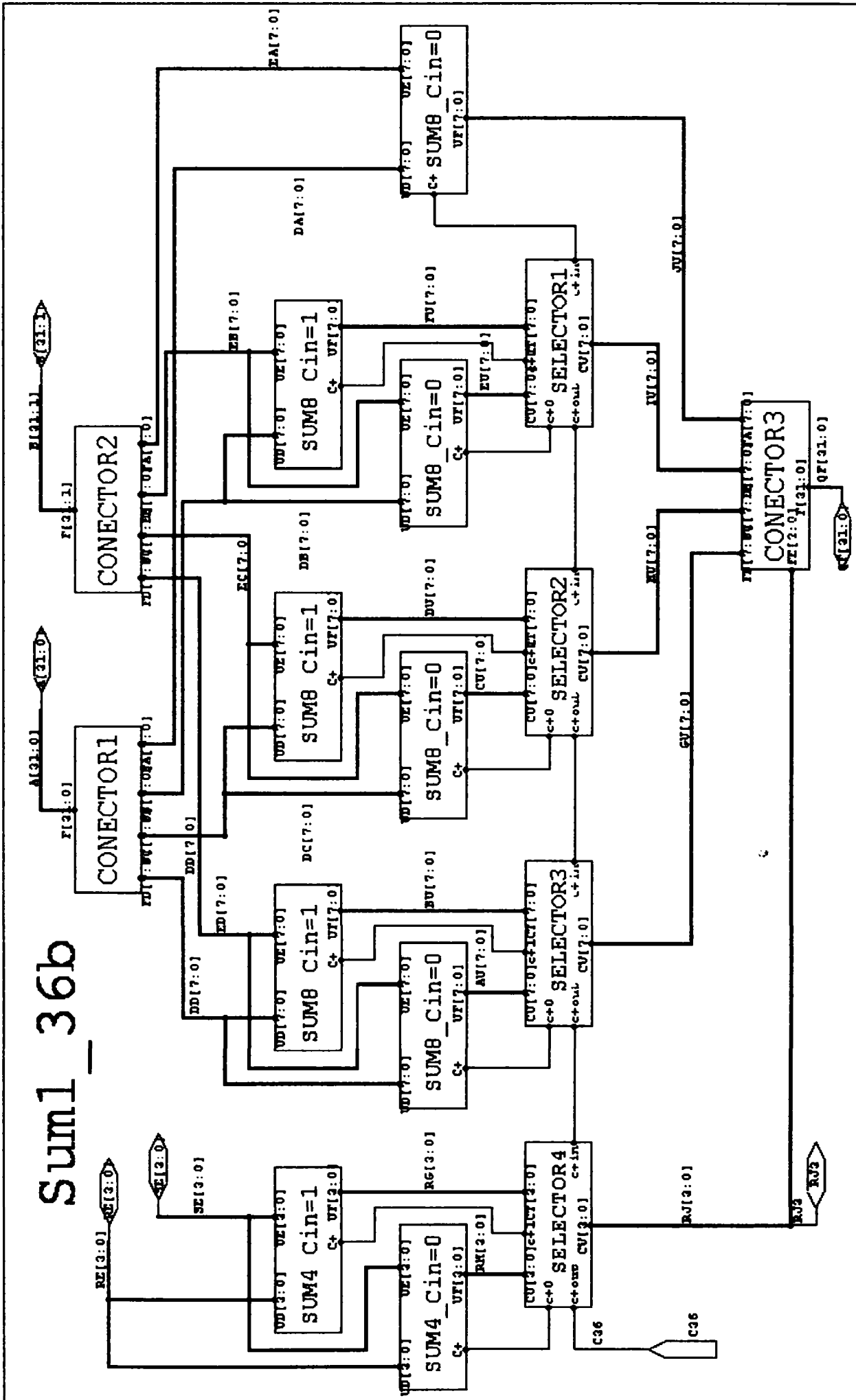


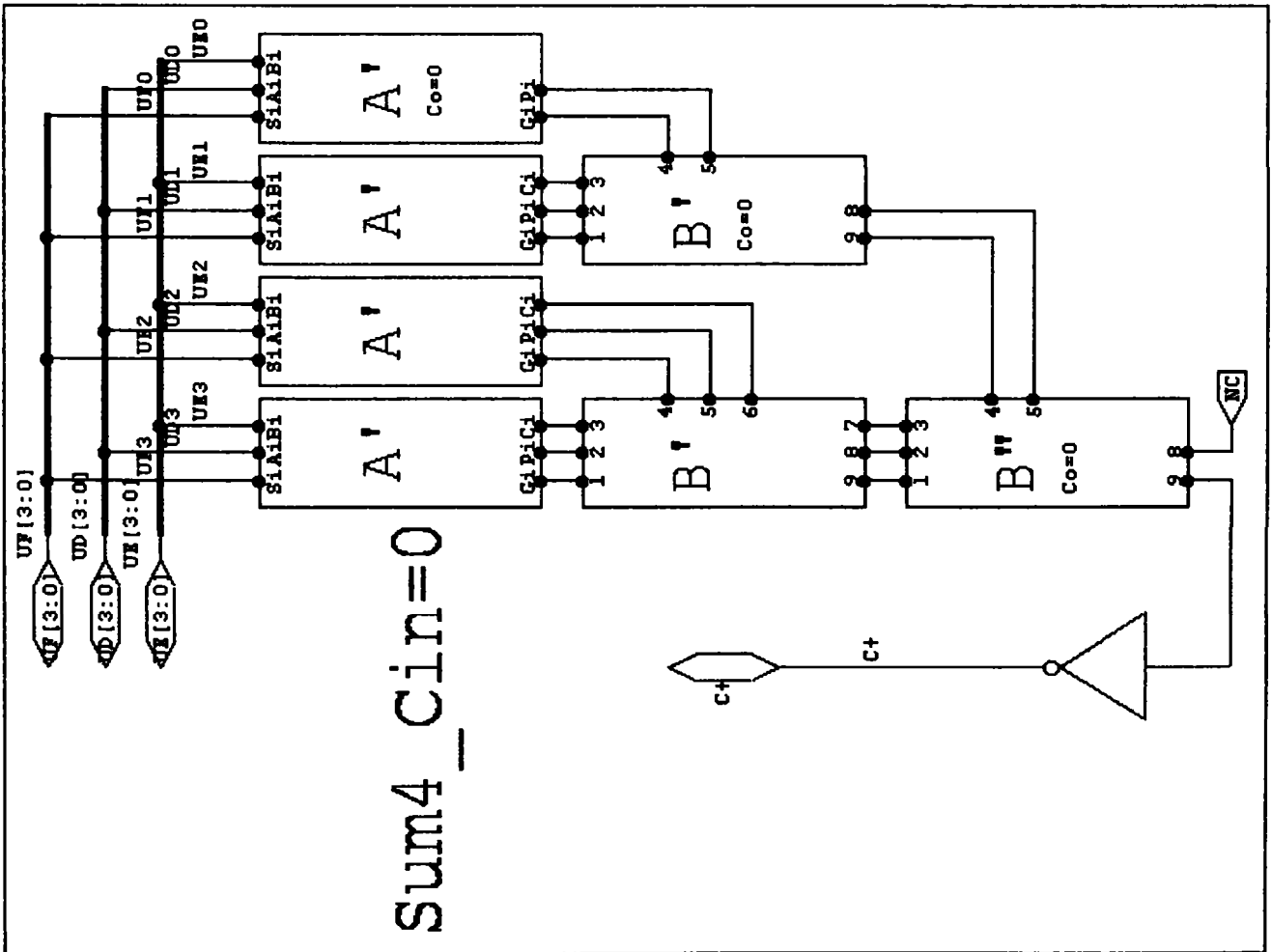
ANEXA 3

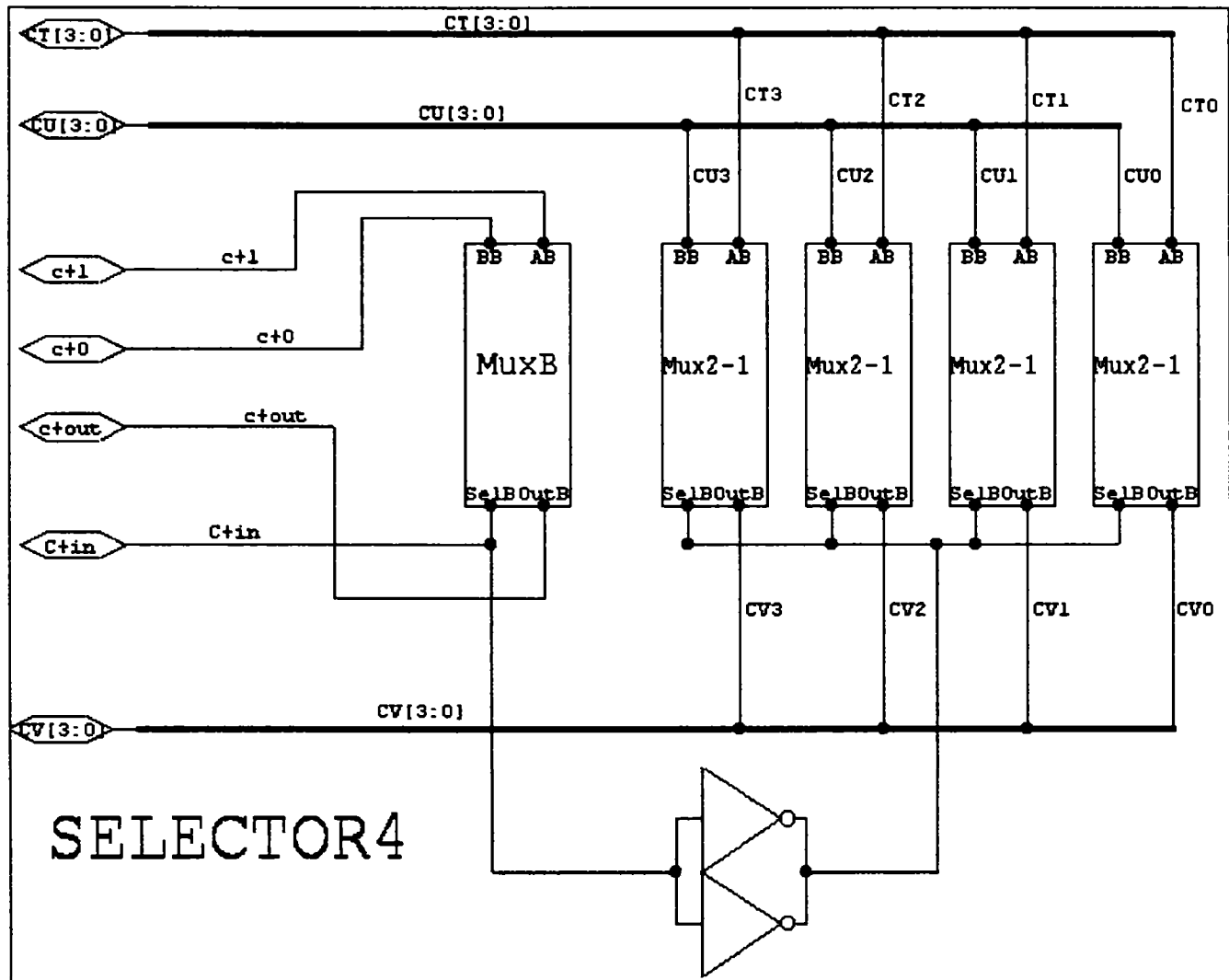


ANEXA 3







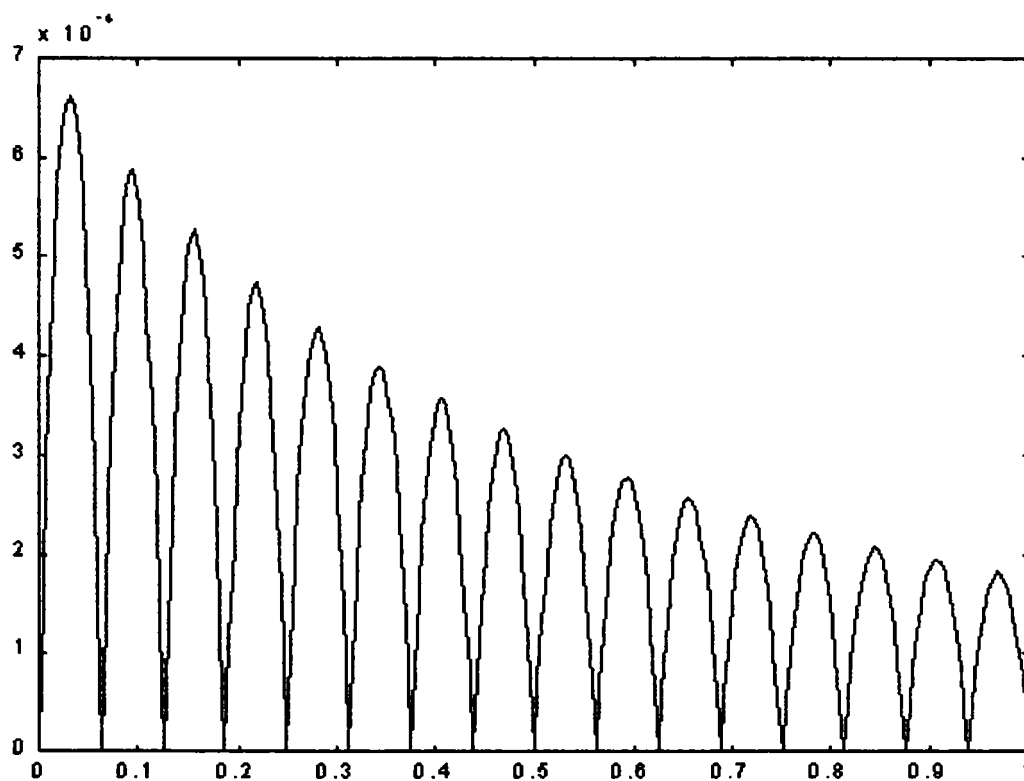


Următoarele blocuri sunt identice cu cele prezentate în Anexa 1:

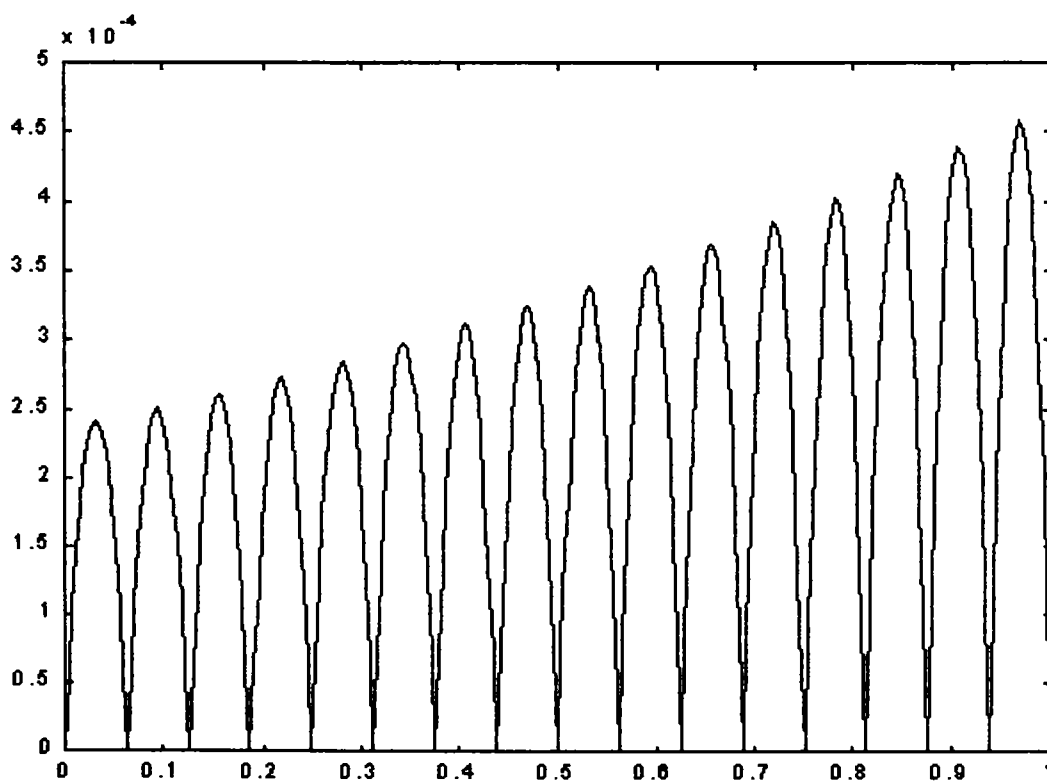
Inversor2;
 Selector;
 SELECTOR1,2,3;
 MUX;
 MuxB;
 Mux2-1;
 SUM8_Cin=0;
 SUM8_Cin=1;
 A', B', B'';
 A' Co=0, A' Co=1;
 B' Co=0, B' Co=1;
 B'' Co=0, B'' Co=1.

ANEXA 4

Eroarea de conversie la logaritmare, obținută în cazul simplificat când interpolarea liniară se face între 16 valori memorate ale funcției $f(y)=\log_2(1+y)-y$, pentru 256 de puncte intermediare:



Eroarea de conversie la antilogaritmare, obținută în cazul simplificat când interpolarea liniară se face între 16 valori memorate ale funcției $f(y)=(1+y)-2^y$, pentru 256 de puncte intermediare:

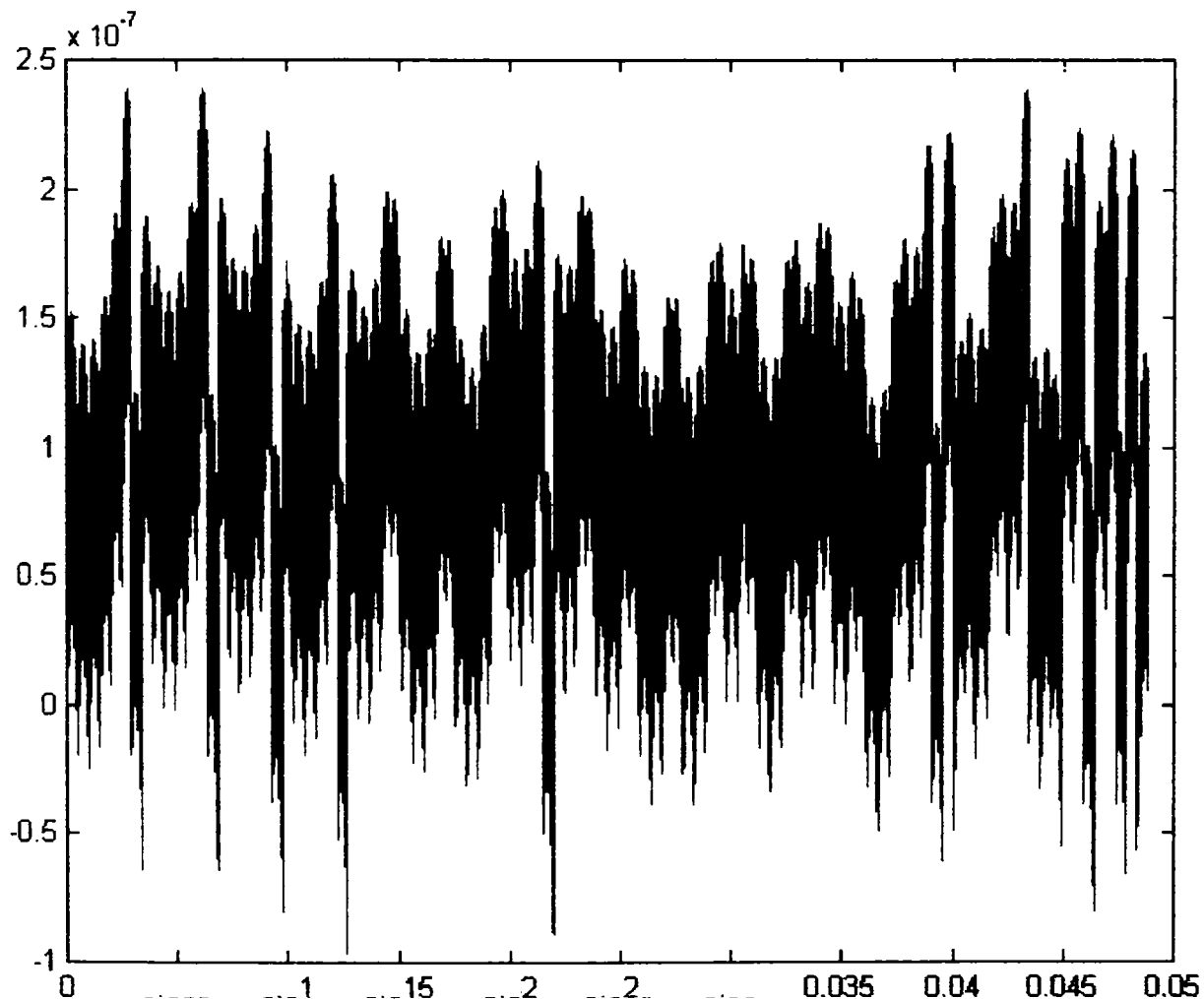


Secvențele Matlab utilizate în cazurile anterioare:

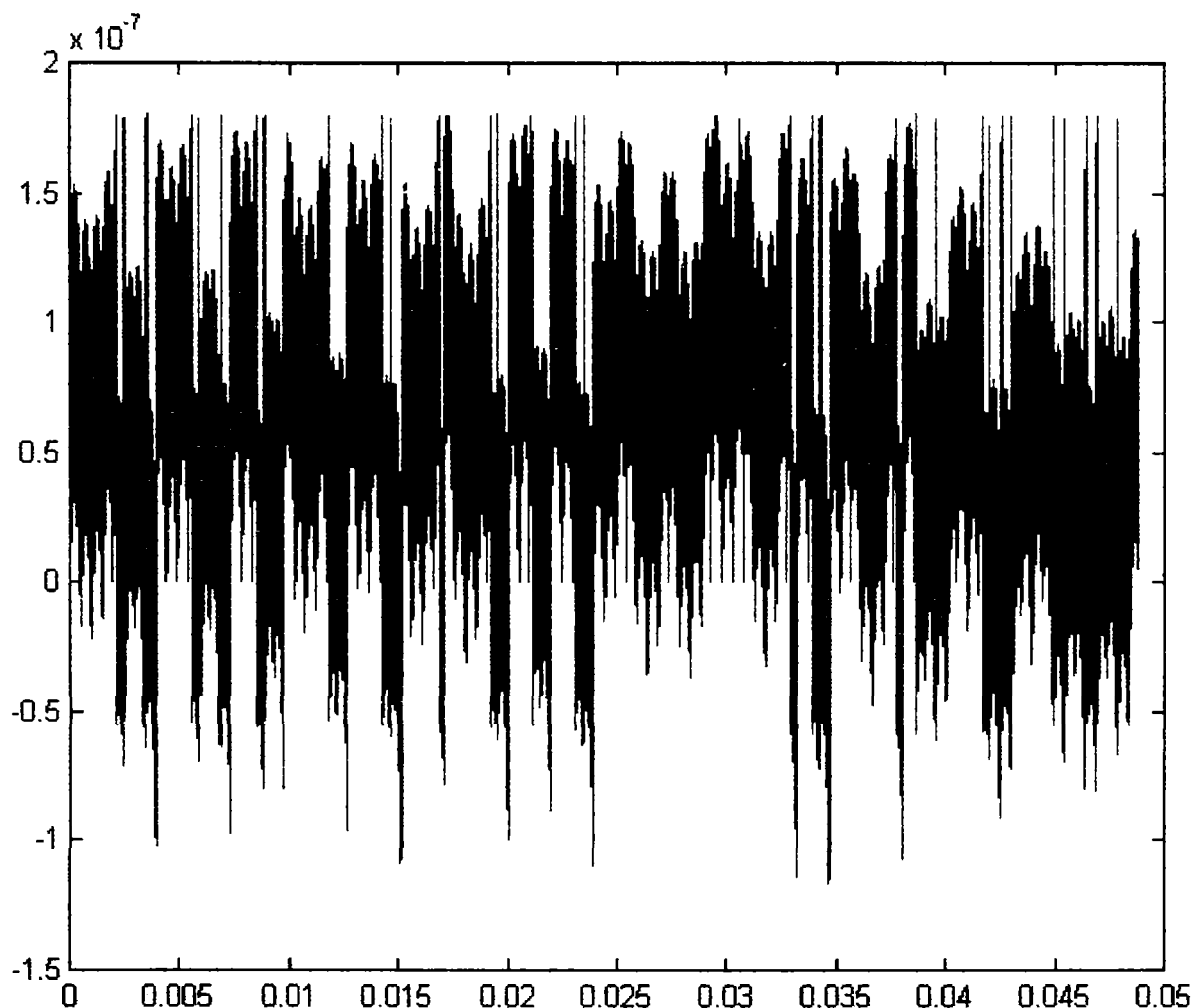
logaritmare	antilogaritmare
<pre> y=0:16:256; L=log2(1+y/256)-y/256; yy=0:1:256; T(:,1)=y'; T(:,2)=L'; z=table1(T,yy); LL=log2(1+yy/256)-yy/256; er=LL-z'; yyy=yy/256; plot(yyy,er) </pre>	<pre> y=0:16:256; AL=1+y/256-2 .^(y/256); yy=0:1:256; T(:,1)=y'; T(:,2)=AL'; z=table1(T,yy); ALL=1+yy/256-2 .^(yy/256); er=ALL-z'; yyy=yy/256; plot(yyy,er) </pre>

Simulările au fost făcute pentru cazul ideal, când reprezentarea binară se face pe un număr infinit de biți, pentru a depista intervalul pe care se obține eroare maximă prin algoritmul de interpolare.

În continuare se prezintă distribuția erorilor de la logaritmare pentru primele 409600 valori calculate prin algoritmul de interpolare propus de F. Lai. Calculul s-a efectuat pe baza primelor 100 de valori (din totalul celor 2048) memorate în memoriile ROMA și ROMA' cu rezoluție de 23 de biți.



Următoarea diagramă prezintă distribuția erorilor de la logaritmare pentru primele 409600 valori calculate prin algoritmul de interpolare îmbunătățit, propus de autorul tezei. Calculul s-a efectuat pe baza primelor 100 de valori (din totalul celor 2048) memorate în memoriile ROMA și ROMA' cu rezoluție de 23 de biți și care includ în anumite locații de memorie valori de corecție pentru limitarea erorii maxime la $1,8 \times 10^{-7}$.



Programul matlab pentru generarea valorilor E_y (ce includ corecții pe anumite intervale, pentru a limita eroarea de conversie la $1,8 \times 10^{-7}$) și ΔE_y ce se memorează în primele 100 de locații ale memoriilor ROMA și ROMA' este prezentat mai jos:

```

y=0:99;
L=log2(1+y/2048)-y/2048;
La=round(L*2^23)/2^23;
x=1:100;
LD=log2(1+x/2048)-x/2048;
D=LD-L;
Da=round(D*2^23)/2^23;
Lteor(1,:)=La;

```

ANEXA 4

```

er(1,1:100)=0;
z=0:409599;
Lz=log2(1+z/8388608)-z/8388608;
Lz=reshape(Lz,4096,100);
for i=1:100,
for j=1:4095,
Lteor(1+j,i)=La(i)+(j/4096)*Da(i);
Lteora(1+j,i)=fix(Lteor(1+j,i)*2^23)/2^23;
er(1+j,i)=Lz(1+j,i)-Lteora(1+j,i);
if er(1+j,i)>1.8*10^-7
La(i)=La(i)+1.2*10^-7;
Lteor(1+j,i)=La(i)+(j/4096)*Da(i);
Lteora(1+j,i)=fix(Lteor(1+j,i)*2^23)/2^23;
er(1+j,i)=Lz(1+j,i)-Lteora(1+j,i);
end
end
end
er=reshape(er,1,409600);
z=z/8388608;
plot(z,er)
La=round(La*2^23);
Lmem=dec2hex(La)
Da=Da*2^23;
Dmem=dec2hex(Da)

```

Adresa	E _y	Adresa	E _y	Adresa	E _y
00	00000	10	06FE5	20	0DCF3
01	00714	11	076CC	21	0E3AC
02	00E25	12	07DAF	22	0EA63
03	01533	13	0848F	23	0F115
04	01C3F	14	08B6C	24	0F7C6
05	02347	15	09247	25	0FE74
06	02A4C	16	0991F	26	1051F
07	0314F	17	09FF4	27	10BC8
08	0384E	18	0A6C7	28	1126D
09	03F4B	19	0AD96	29	1190F
0A	04645	1A	0B462	2A	11FAF
0B	04D3D	1B	0BB2C	2B	1264D
0C	05431	1C	0C1F3	2C	12CE7
0D	05B22	1D	0C8B8	2D	1337E
0E	06211	1E	0CF79	2E	13A13
0F	068FC	1F	0D637	2F	140A6

ANEXA 4

Adresa	E _y	Adresa	E _y	Adresa	E _y
30	14735	40	1AEB5	50	2137F
31	14DC1	41	1B516	51	219B5
32	1544B	42	1BB74	52	21FE8
33	15AD2	43	1C1D0	53	22618
34	16156	44	1C828	54	22C46
35	167D8	45	1CE7F	55	23272
36	16E57	46	1D4D2	56	2389A
37	174D3	47	1DB22	57	23EC0
38	17B4C	48	1E170	58	244E3
39	181C3	49	1E7BB	59	24B03
3A	18837	4A	1EE04	5A	25121
3B	18EA8	4B	1F44A	5B	2573C
3C	19516	4C	1FA8D	5C	25D55
3D	19B82	4D	200CE	5D	2636B
3E	1A1EB	4E	2070B	5E	2697E
3F	1A851	4F	20D47	5F	26F8F
				60	2759D
				61	27BA8
				62	281B1
				63	287B7

Adresa	ΔE _y	Adresa	ΔE _y	Adresa	ΔE _y
00	714	10	6E6	20	6B9
01	711	11	6E3	21	6B6
02	70E	12	6E0	22	6B3
03	70B	13	6DE	23	6B1
04	708	14	6DB	24	6AE
05	705	15	6D8	25	6AB
06	703	16	6D5	26	6A8
07	700	17	6D2	27	6A5
08	6FD	18	6CF	28	6A3
09	6FA	19	6CD	29	6A0
0A	6F7	1A	6CA	2A	69D
0B	6F4	1B	6C7	2B	69A
0C	6F1	1C	6C4	2C	698
0D	6EF	1D	6C1	2D	695
0E	6EC	1E	6BF	2E	692
0F	6E9	1F	6BC	2F	68F

ANEXA 4

Adresa	ΔE_y		Adresa	ΔE_y		Adresa	ΔE_y
30	68D		40	661		50	636
31	68A		41	65E		51	633
32	687		42	65B		52	630
33	684		43	659		53	62E
34	682		44	656		54	62B
35	67F		45	653		55	628
36	67C		46	651		56	626
37	679		47	64E		57	623
38	677		48	64B		58	620
39	674		49	649		59	61E
3A	671		4A	646		5A	61B
3B	66E		4B	643		5B	619
3C	66C		4C	640		5C	616
3D	669		4D	63E		5D	613
3E	666		4E	63B		5E	611
3F	664		4F	638		5F	60E
						60	60B
						61	609
						62	606
						63	603

Secvența matlab pentru obținerea diagramei de la pagina A4.2:

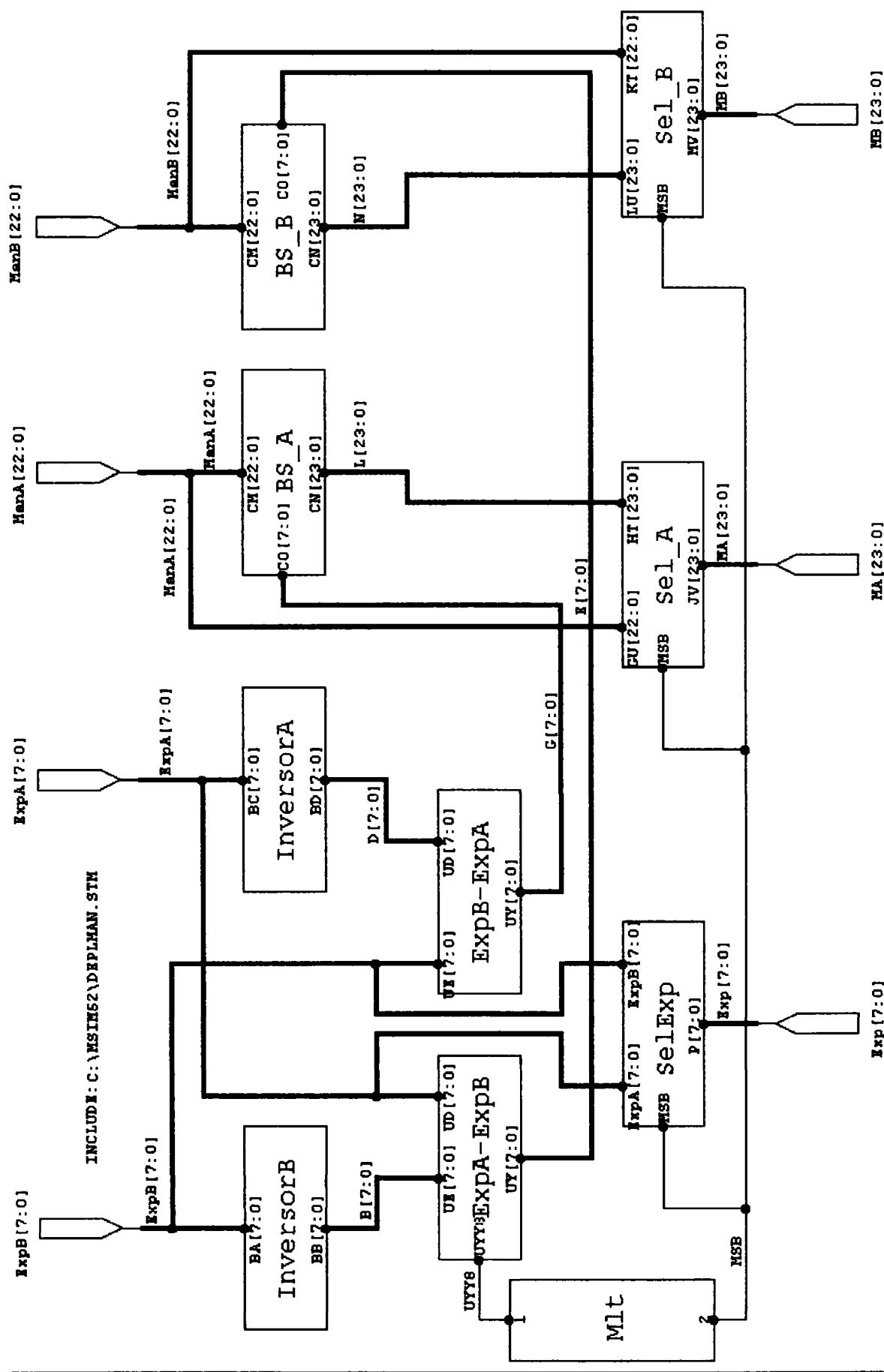
```

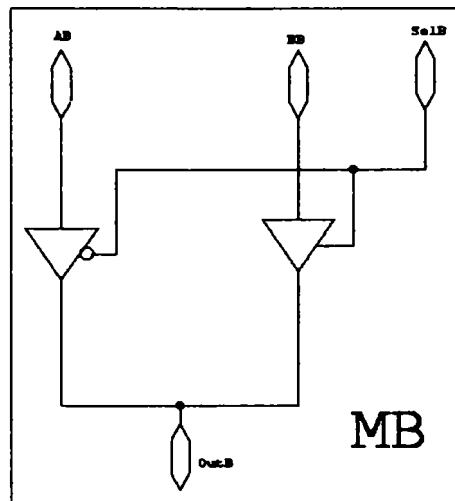
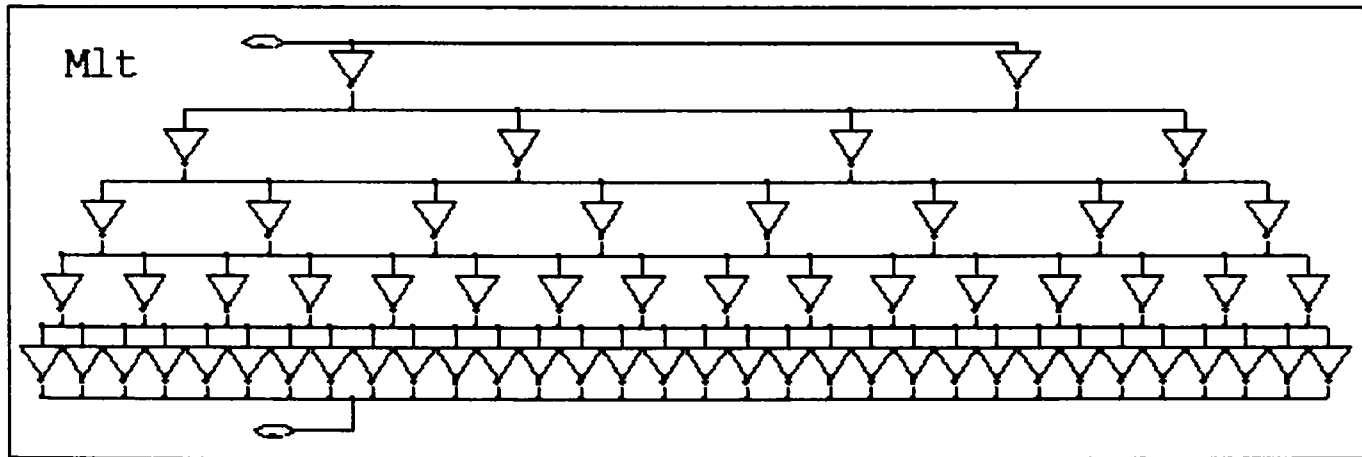
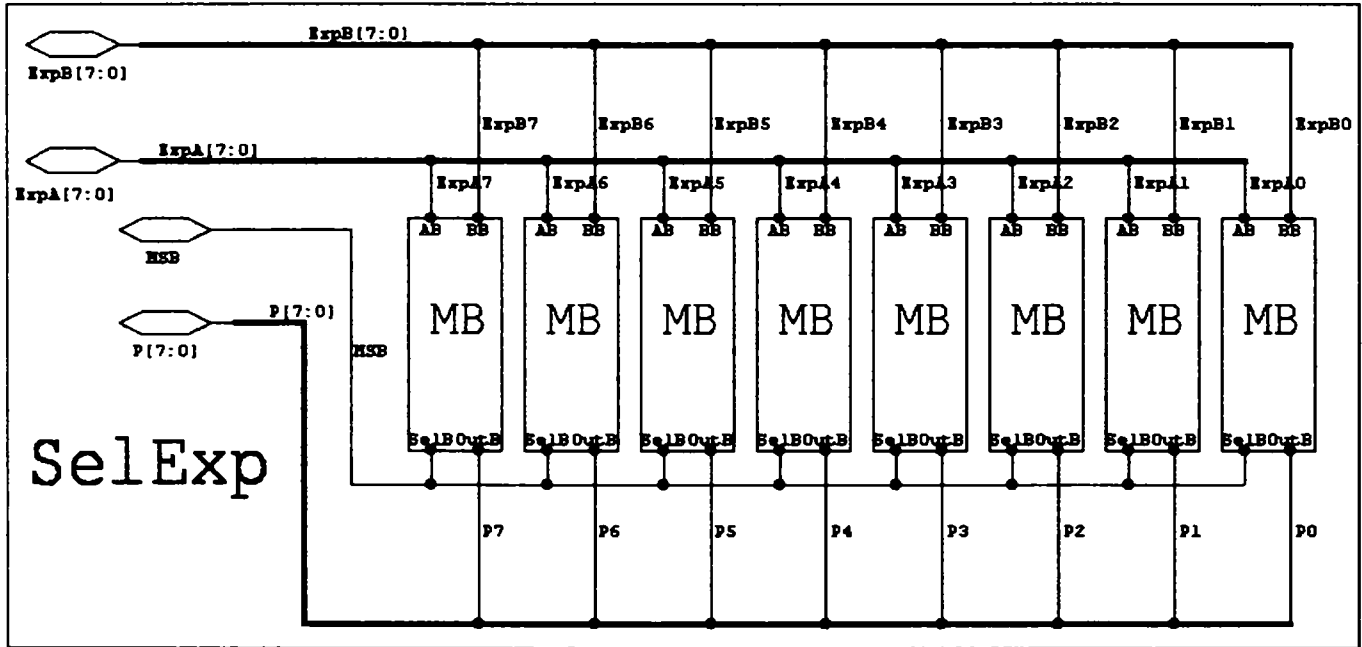
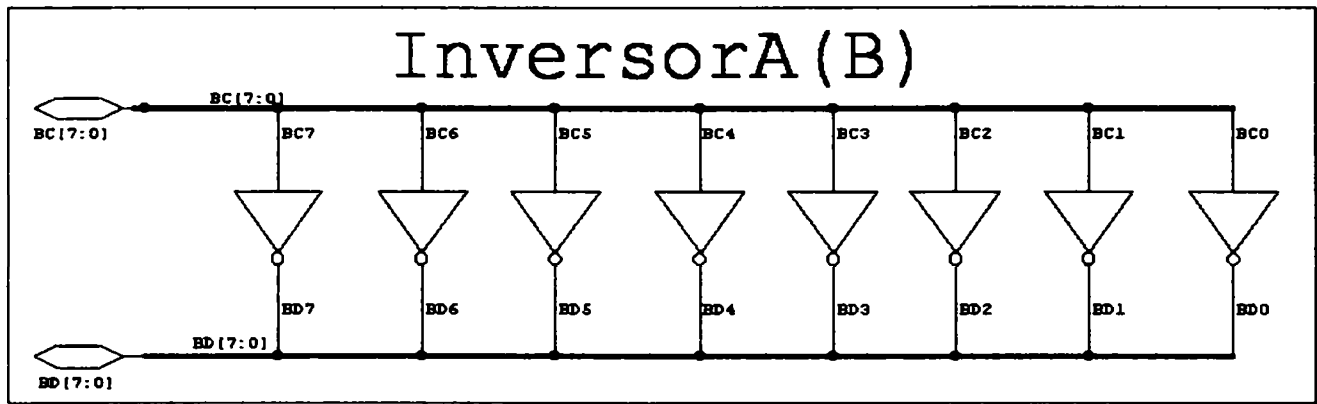
y=0:99;
L=log2(1+y/2048)-y/2048;
La=round(L*2^23)/2^23;
x=1:100;
LD=log2(1+x/2048)-x/2048;
D=LD-L;
Da=round(D*2^23)/2^23;
Lteor(1,:)=La;
for i=1:100,
for j=1:4095,
Lteor(1+j,i)=La(i)+(j/4096)*Da(i);
end
end
Lteora=fix(Lteor*2^23)/2^23;
Lalgoritm=reshape(Lteora,1,409600);
z=0:409599;
Lz=log2(1+z/8388608)-z/8388608;
er=Lz-Lalgoritm;
z=z/8388608;
plot(z,er)

```

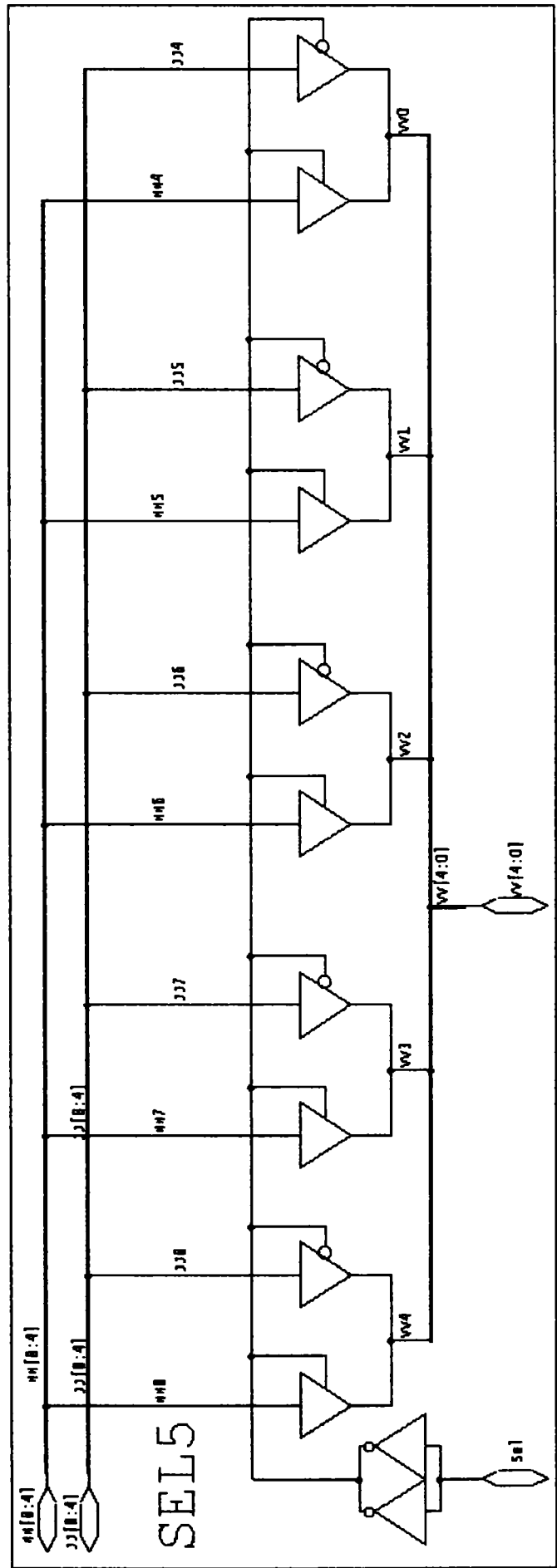
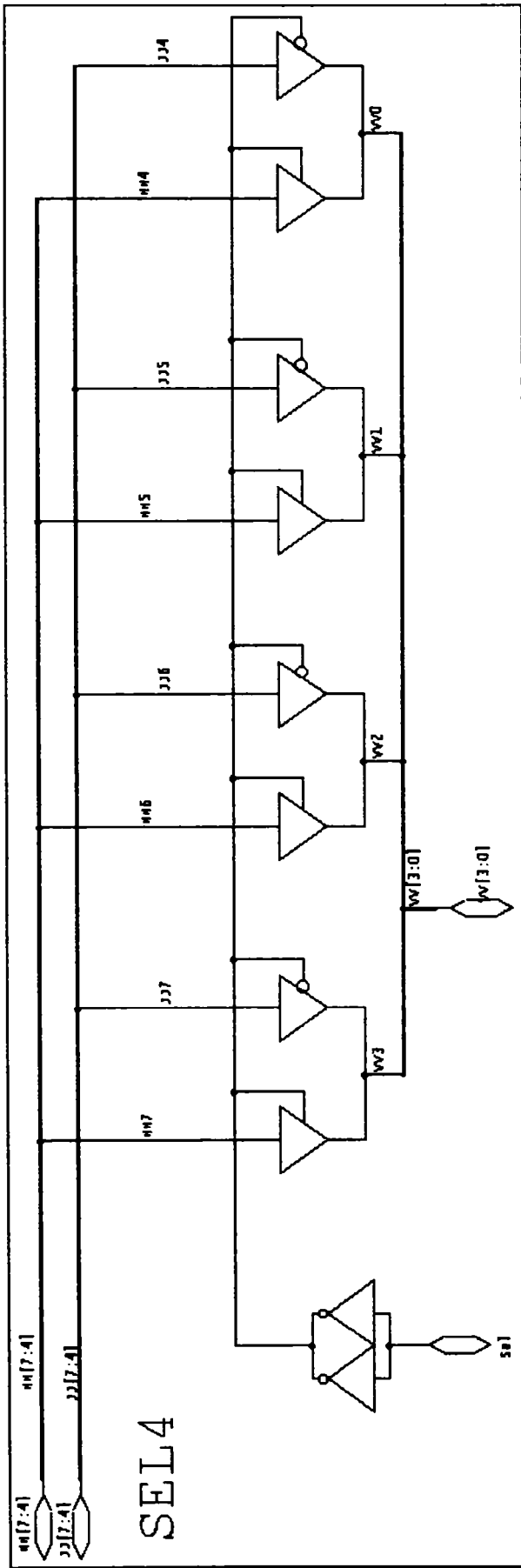

ANEXA 5

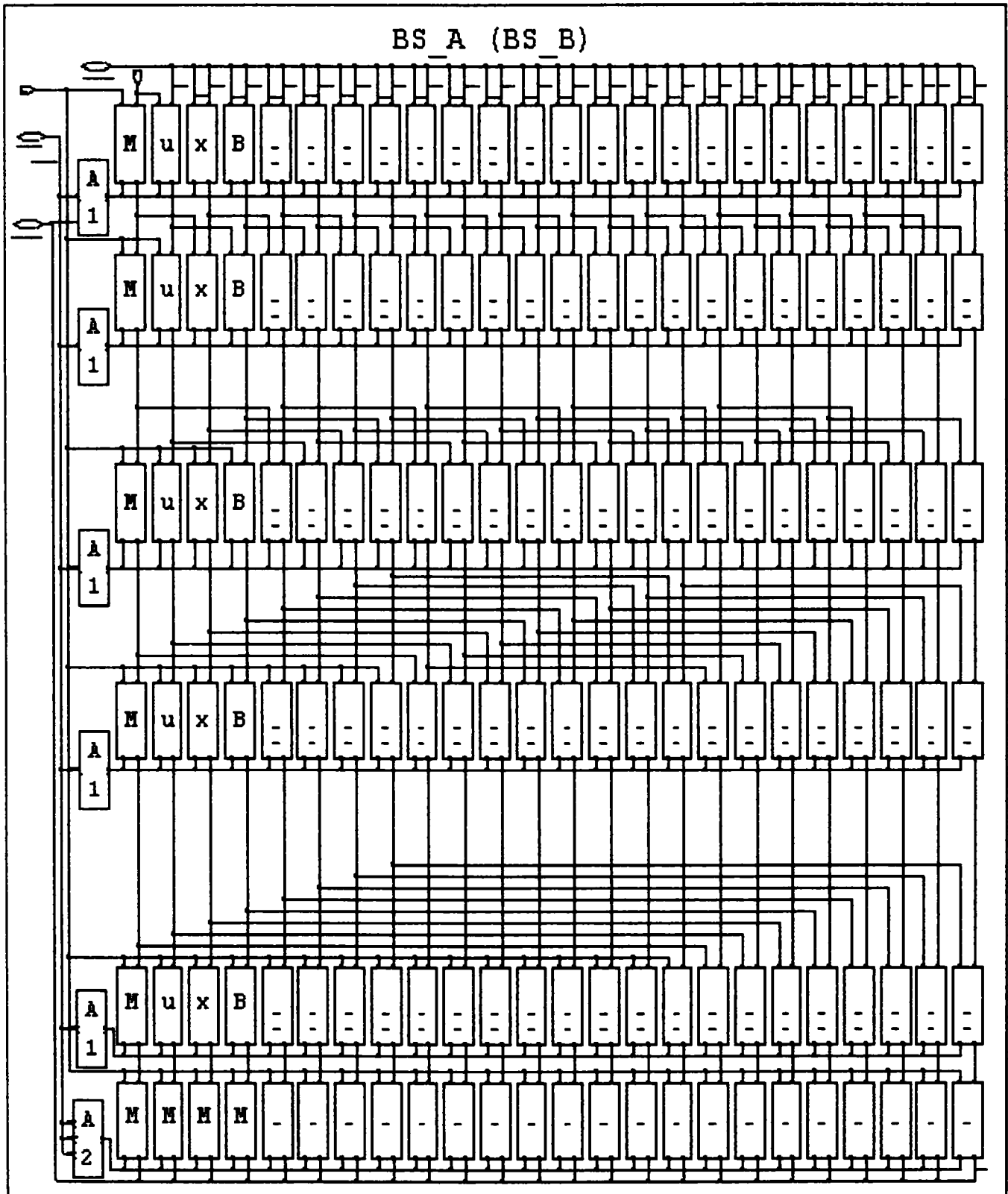
CIRCUIT DEPL. MANT. - NIVEL PIPELINE 1- VARIANTA 1



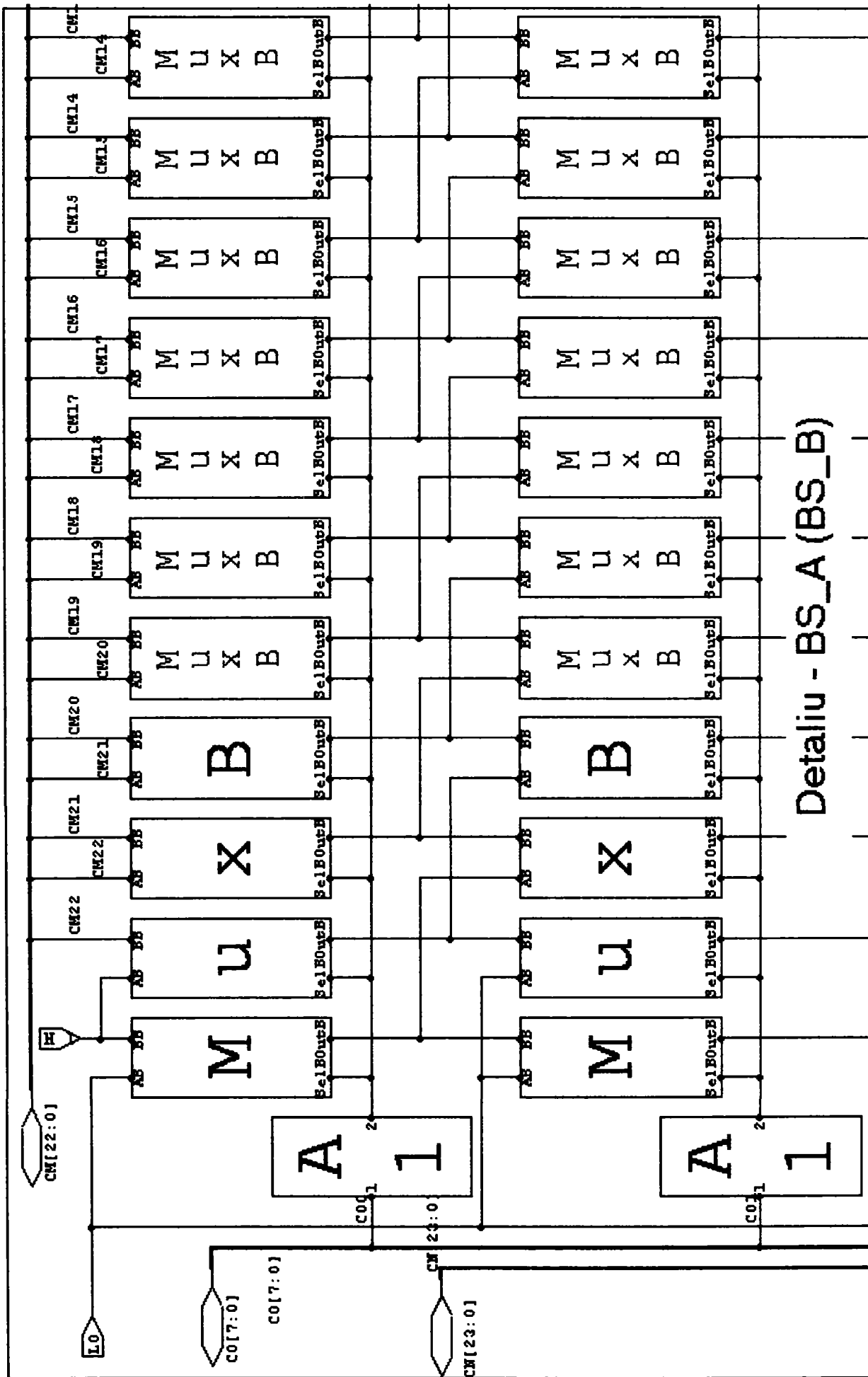


ANEXA 5

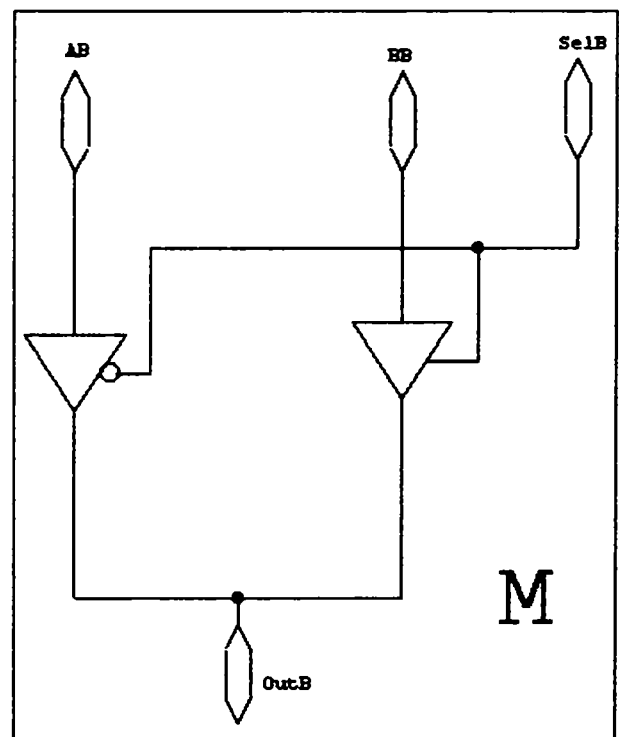
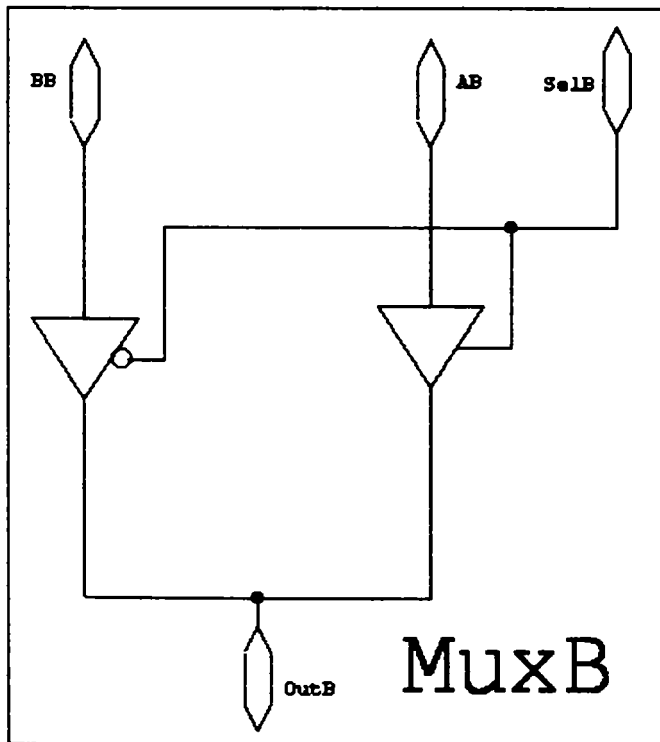
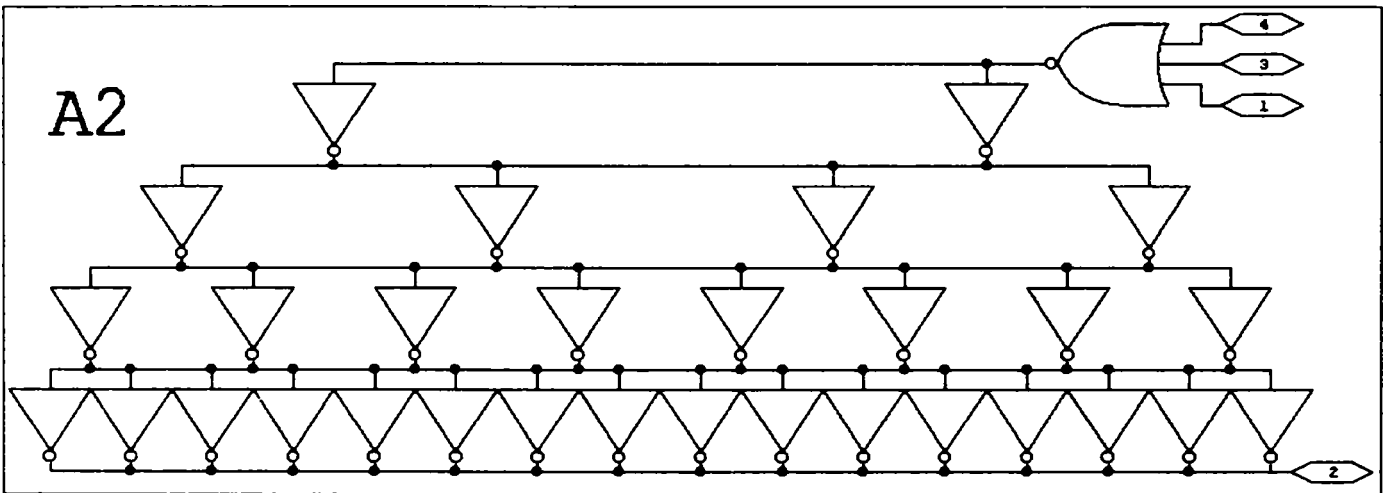
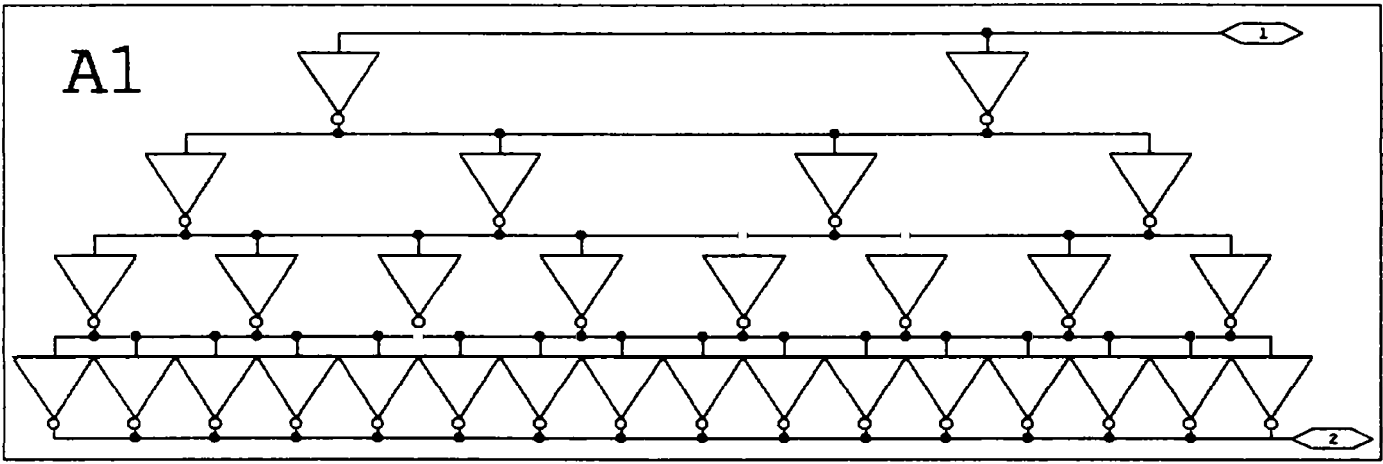




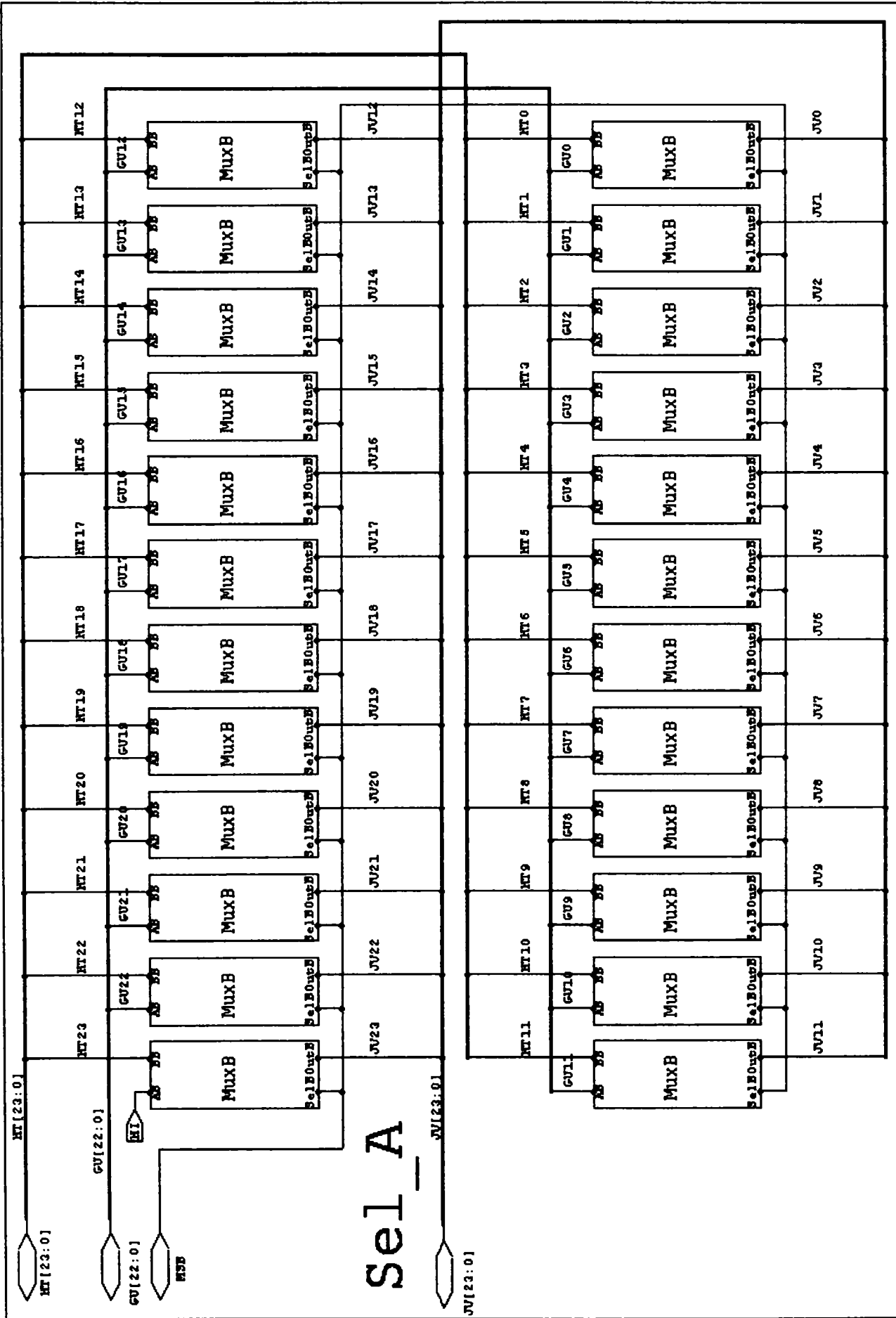
ANEXA 5



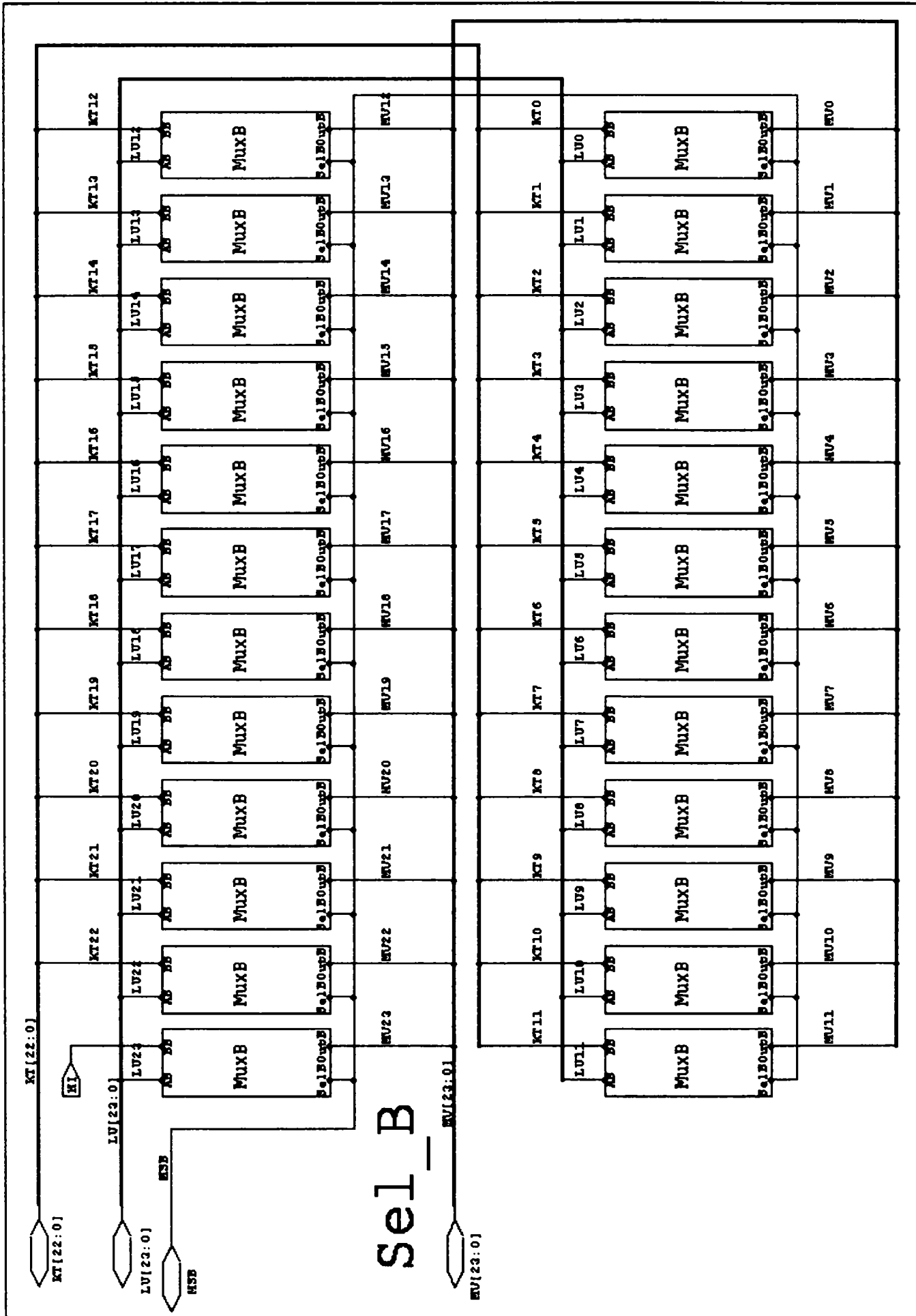
Detaliu - BS_A (BS_B)



ANEXA 5

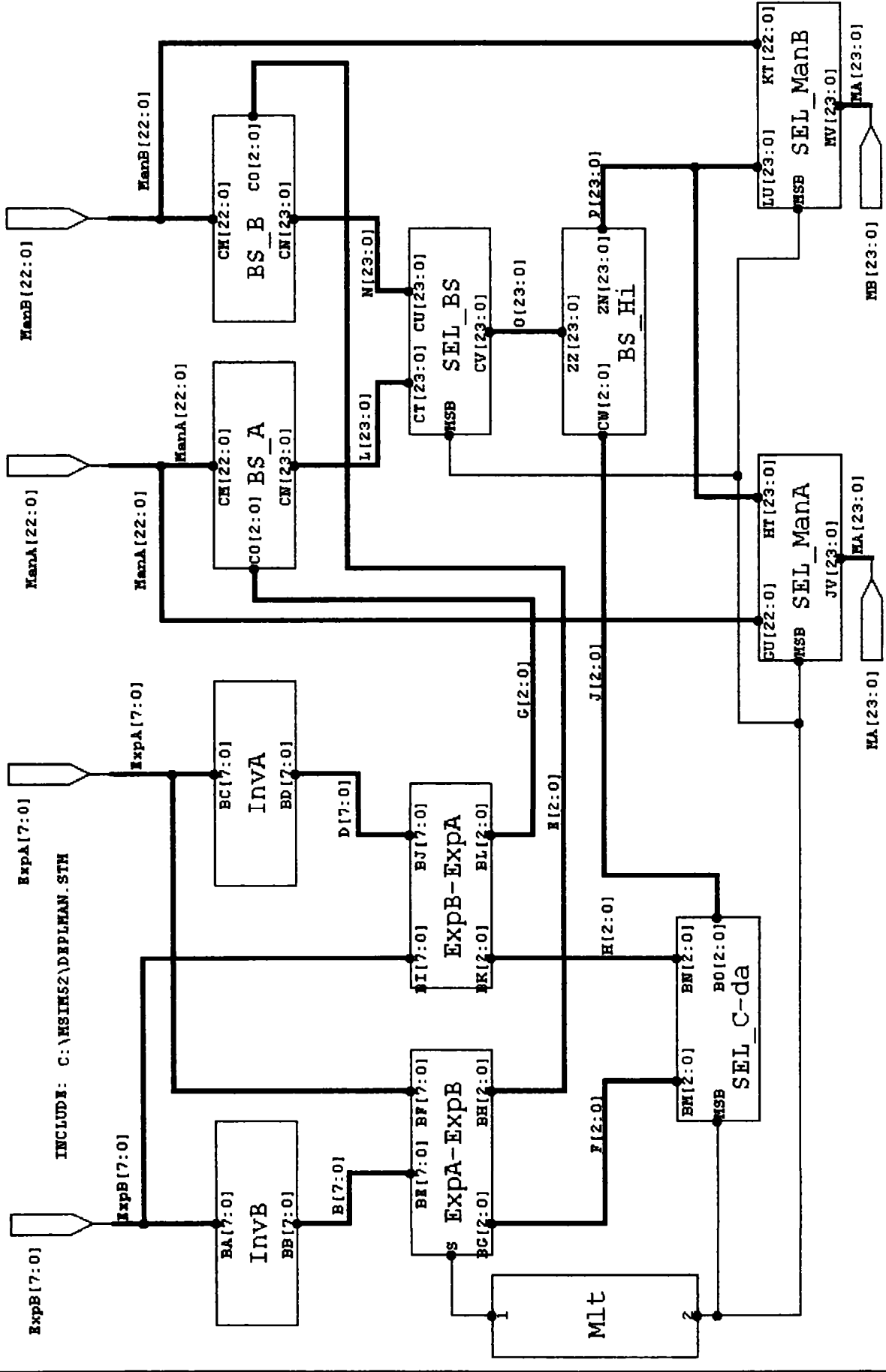


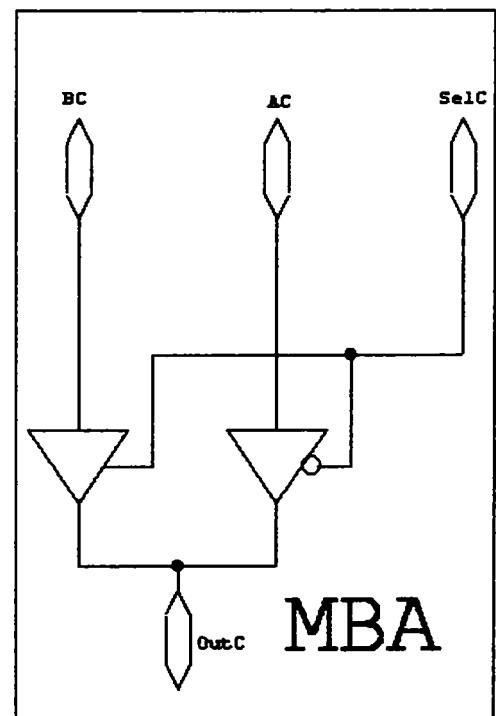
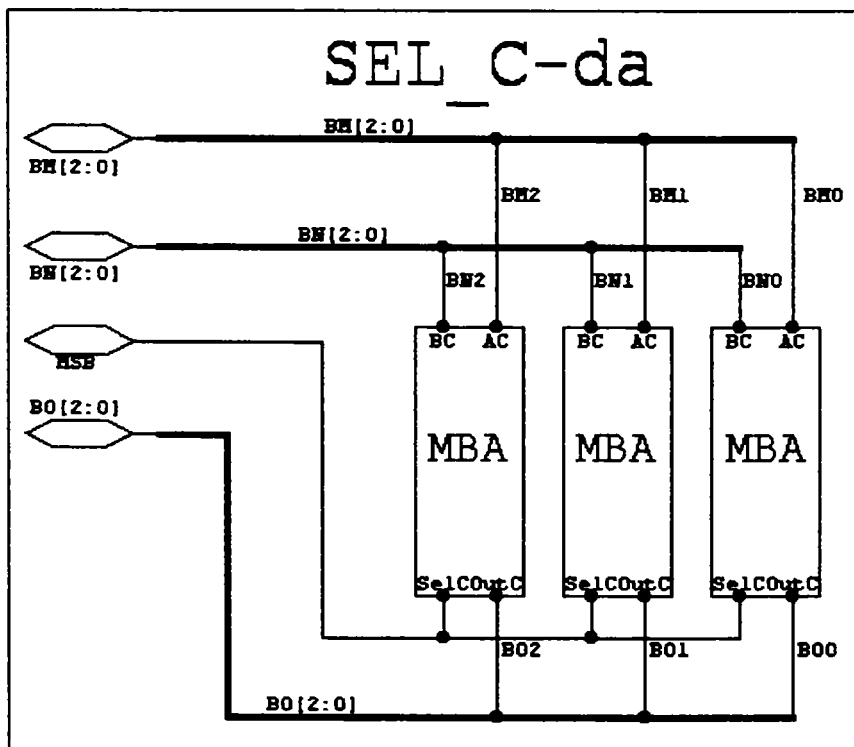
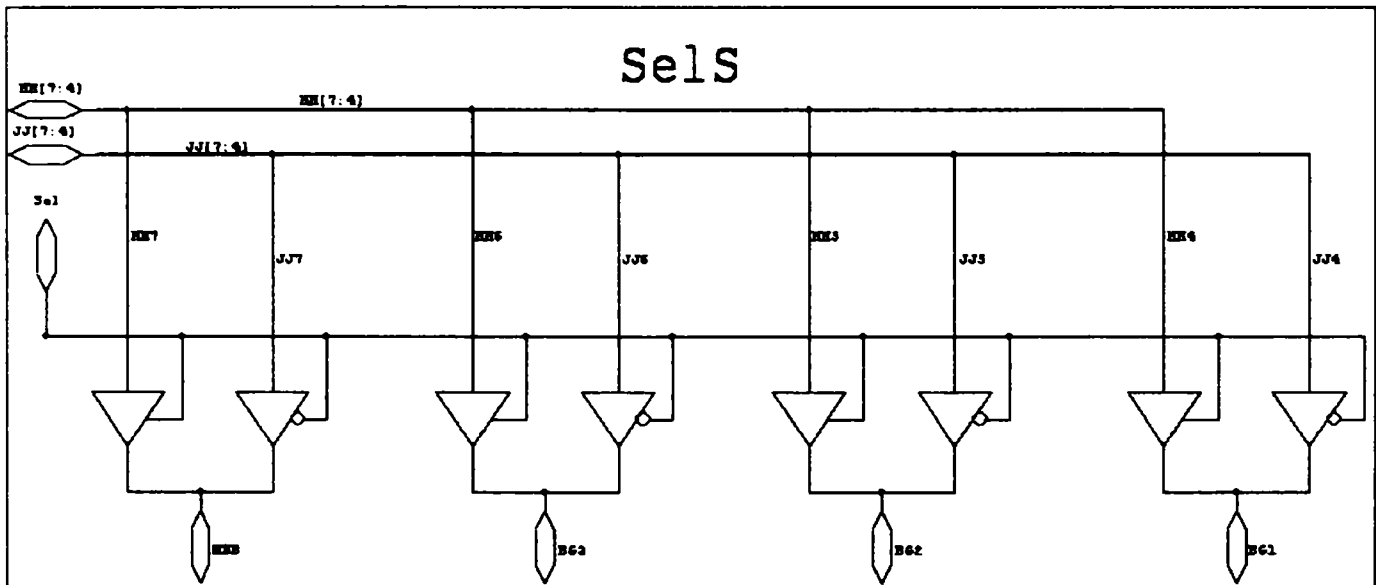
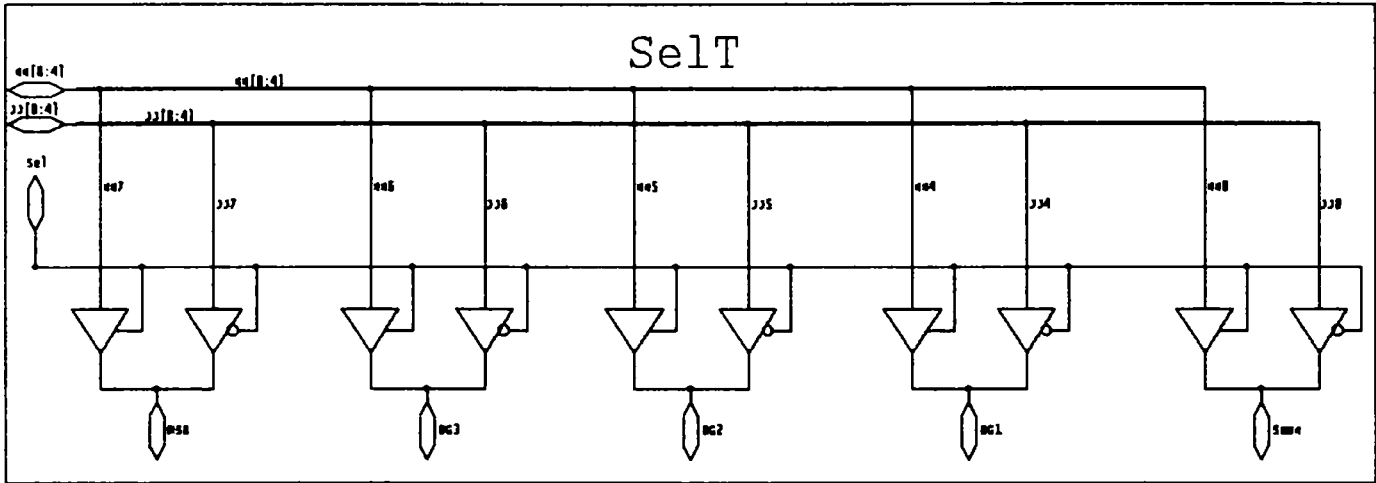
ANEXA 5



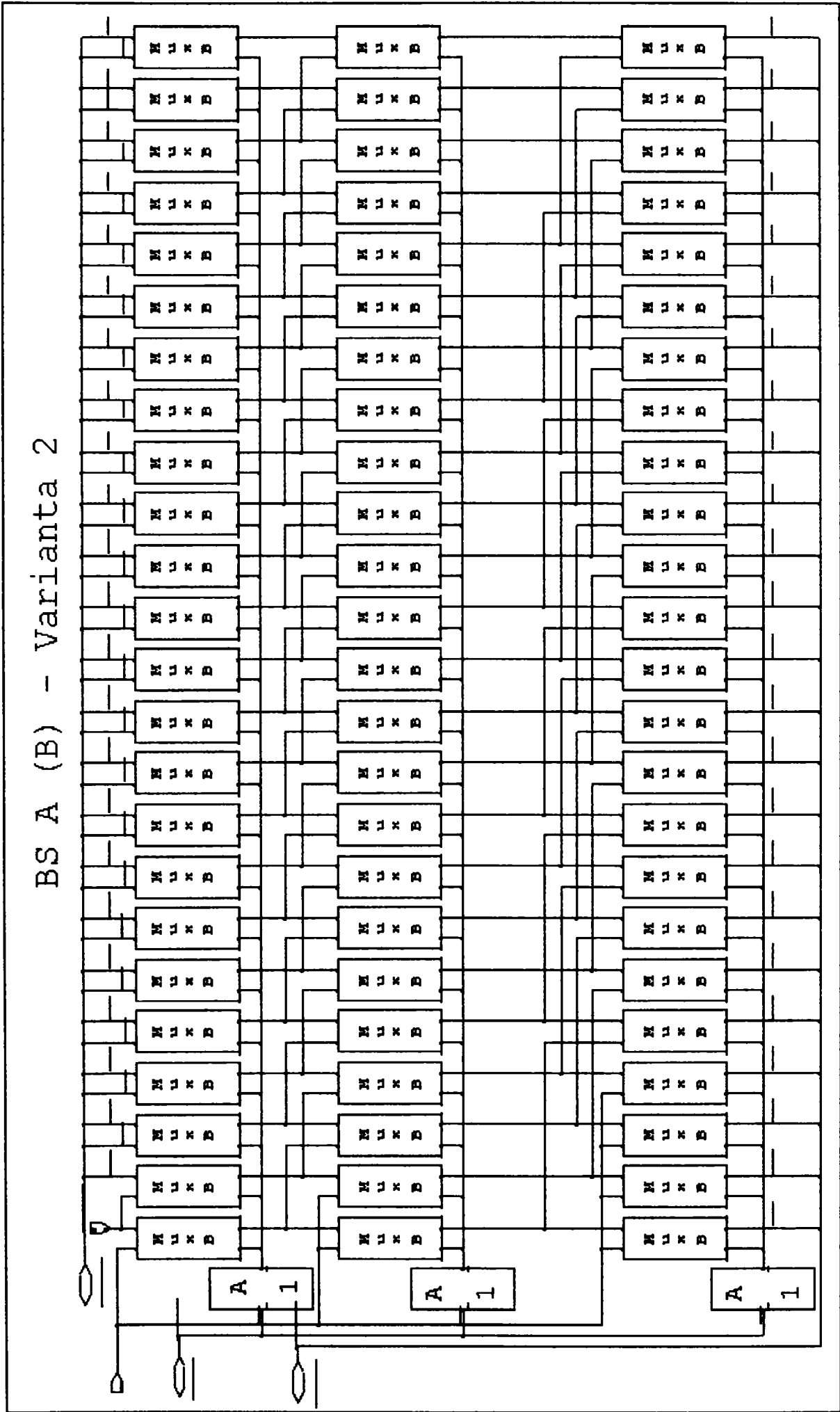
ANEXA 5

CIRCUIT DEPL. MANT. - NIVEL PIPELINE 1- VARIANTA 2

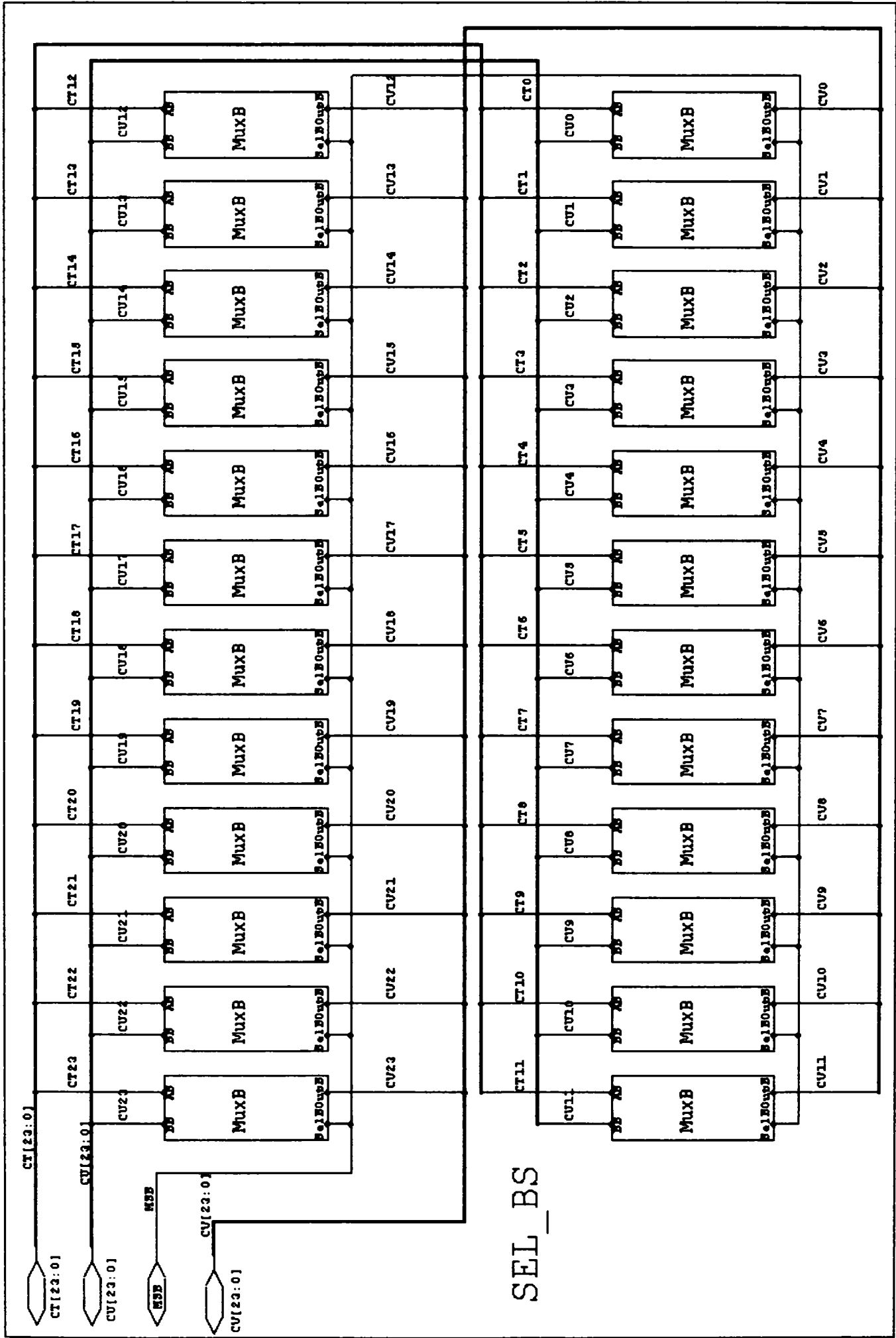




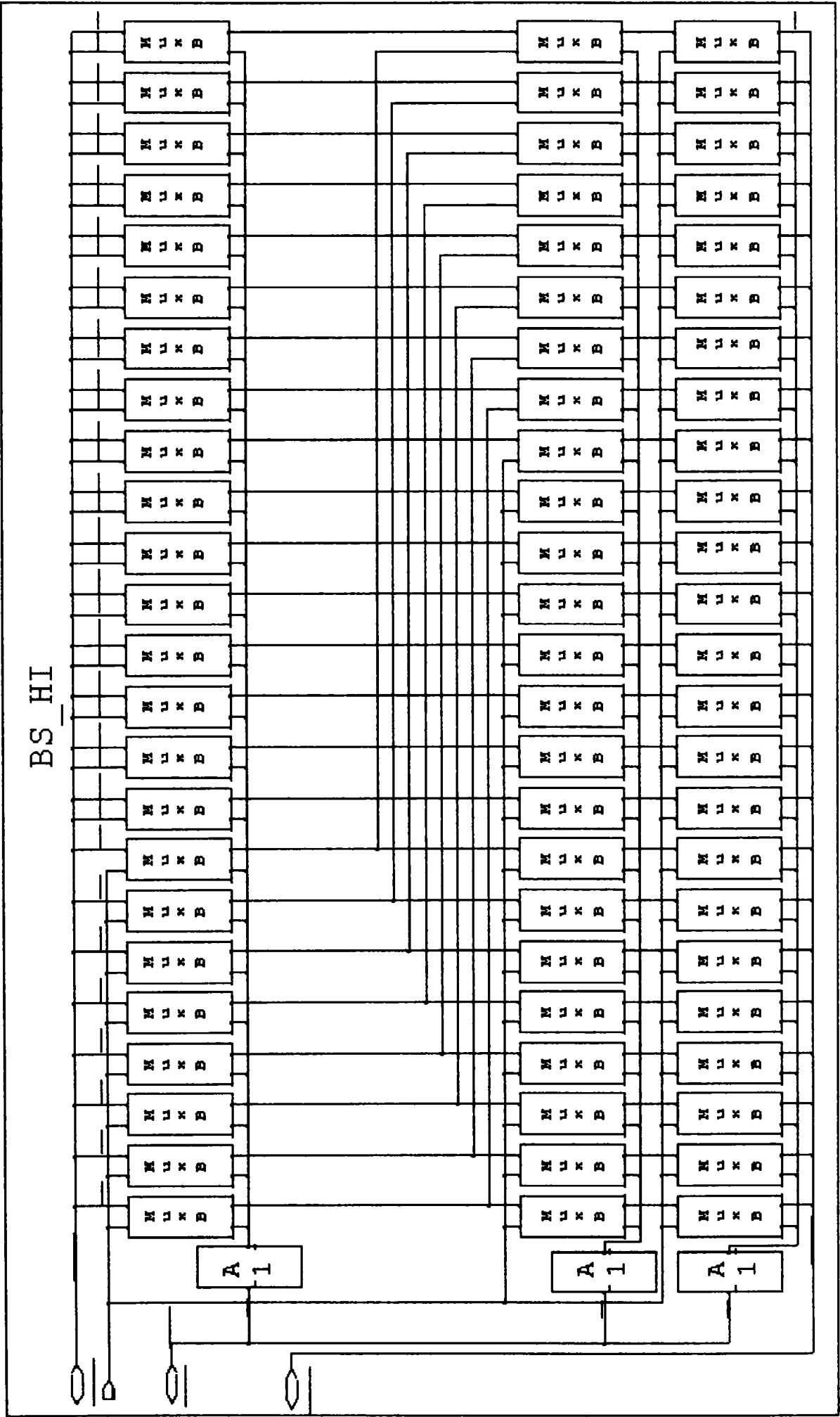
BS A (B) - Varianta 2



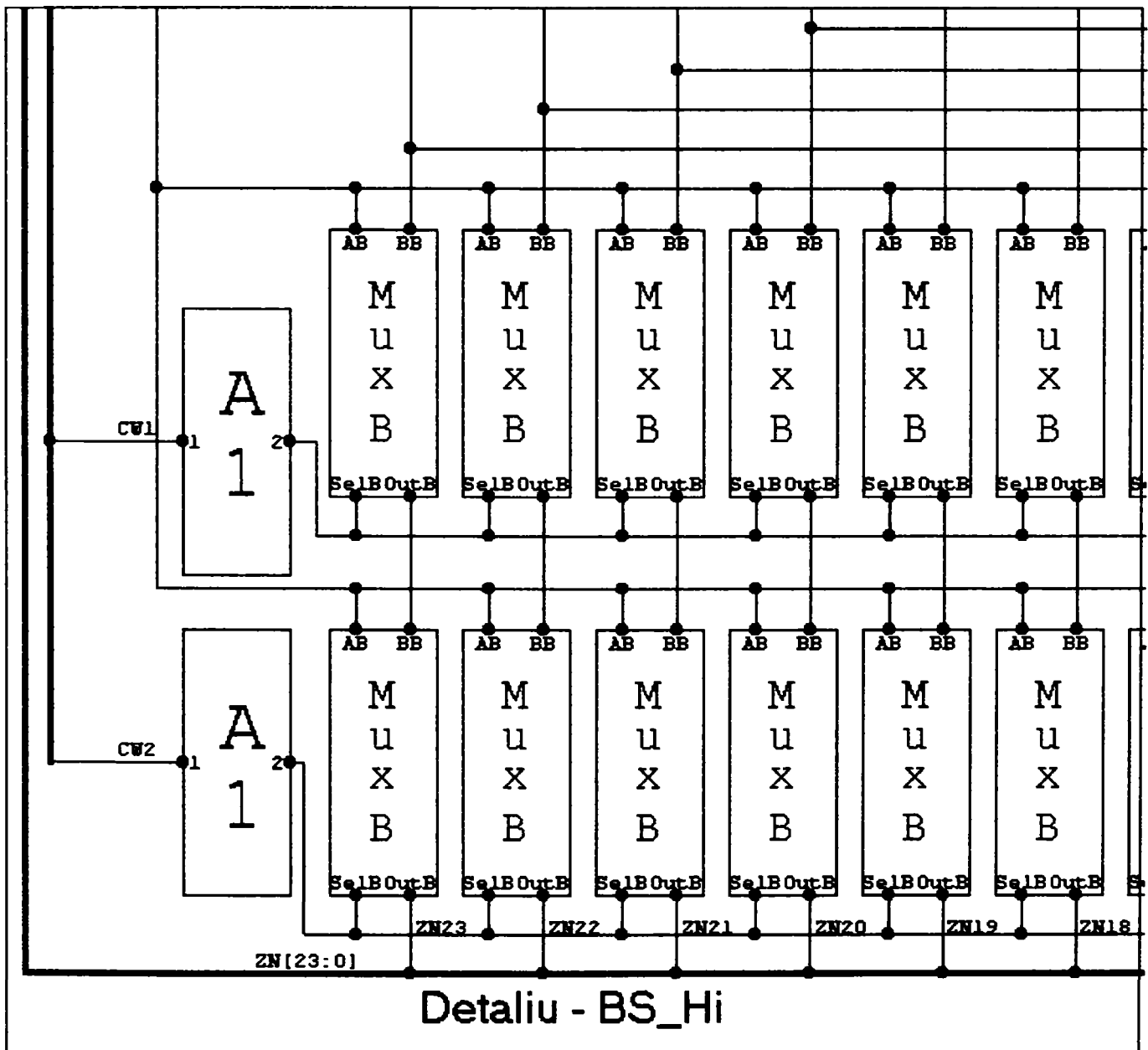
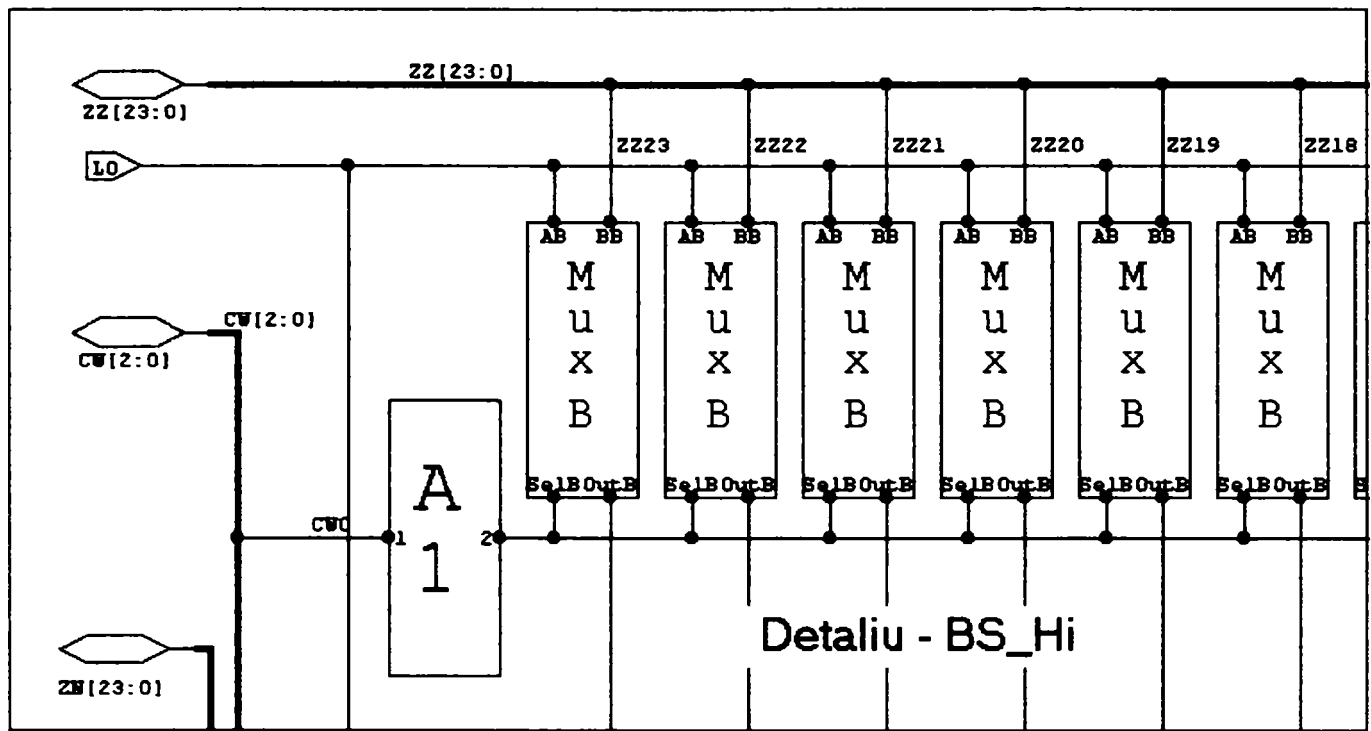
ANEXA 5



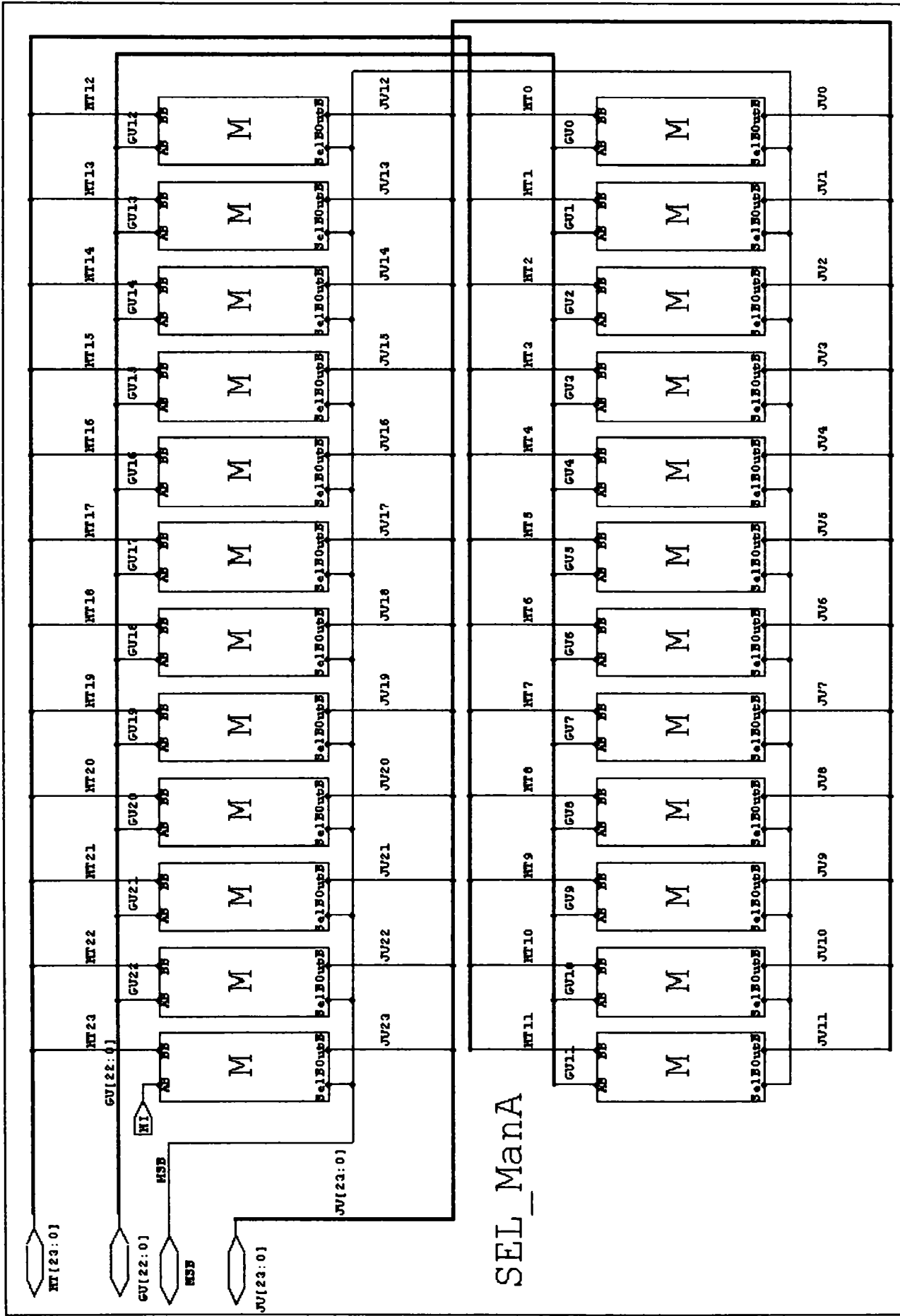
BS HI



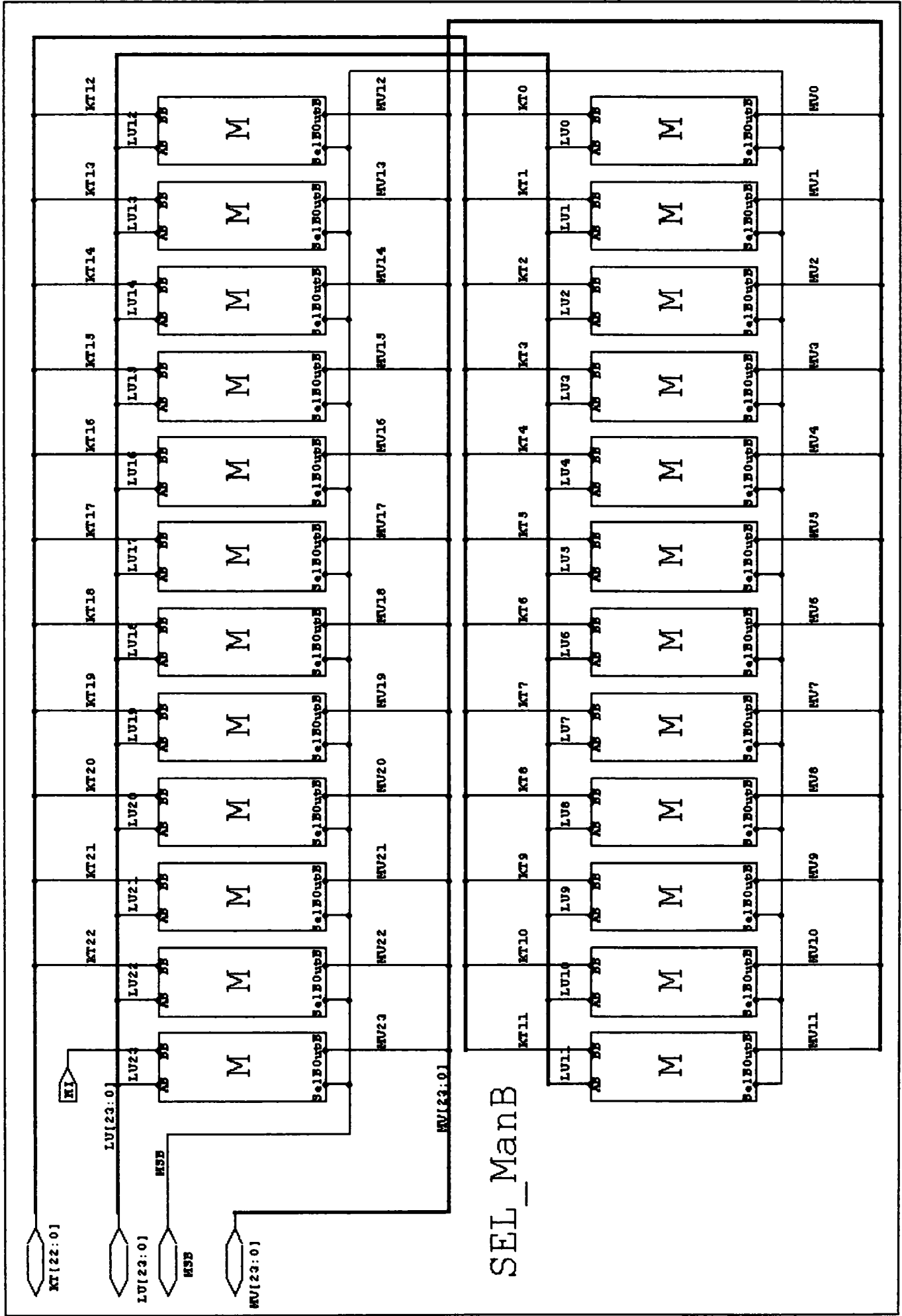
ANEXA 5



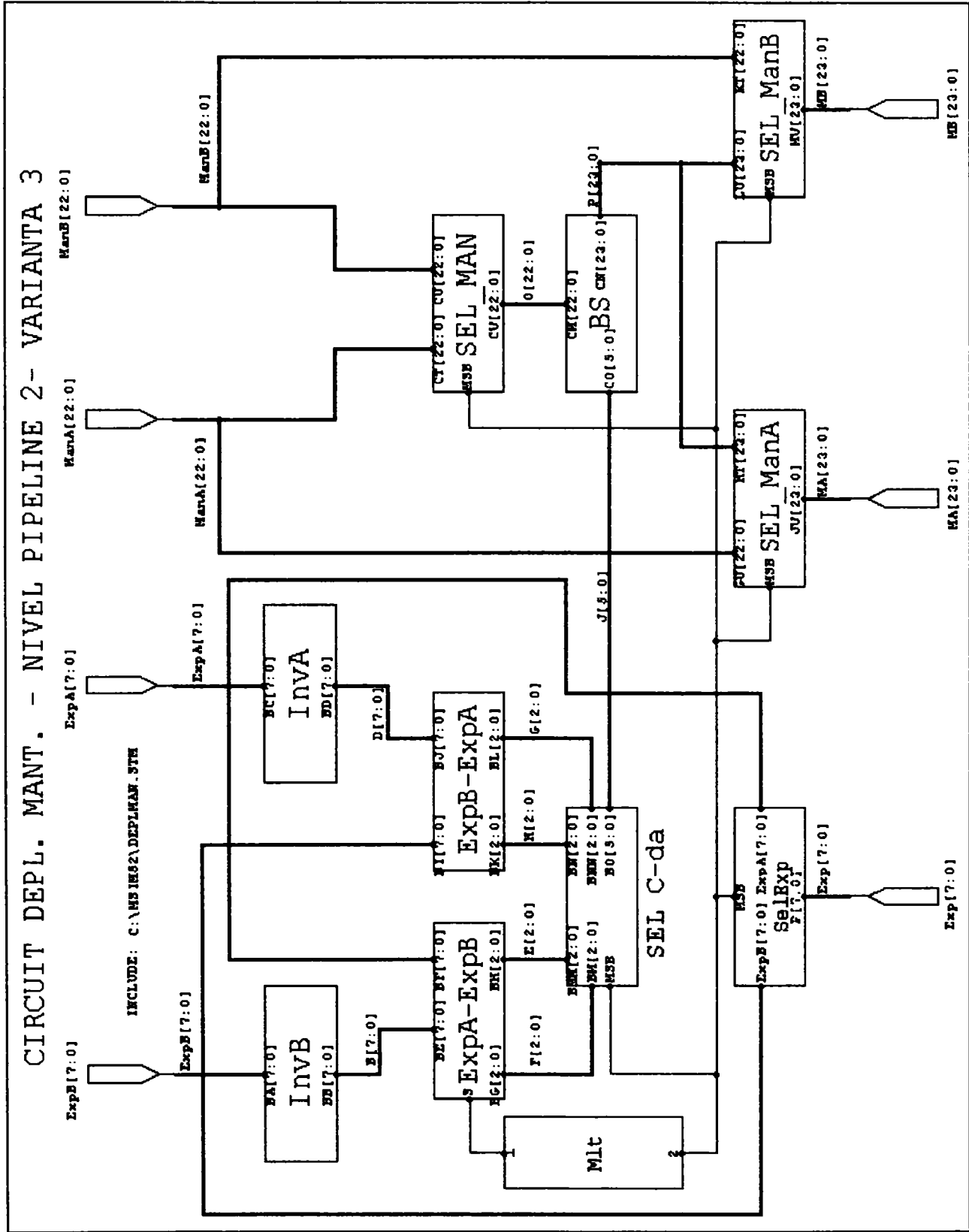
ANEXA 5



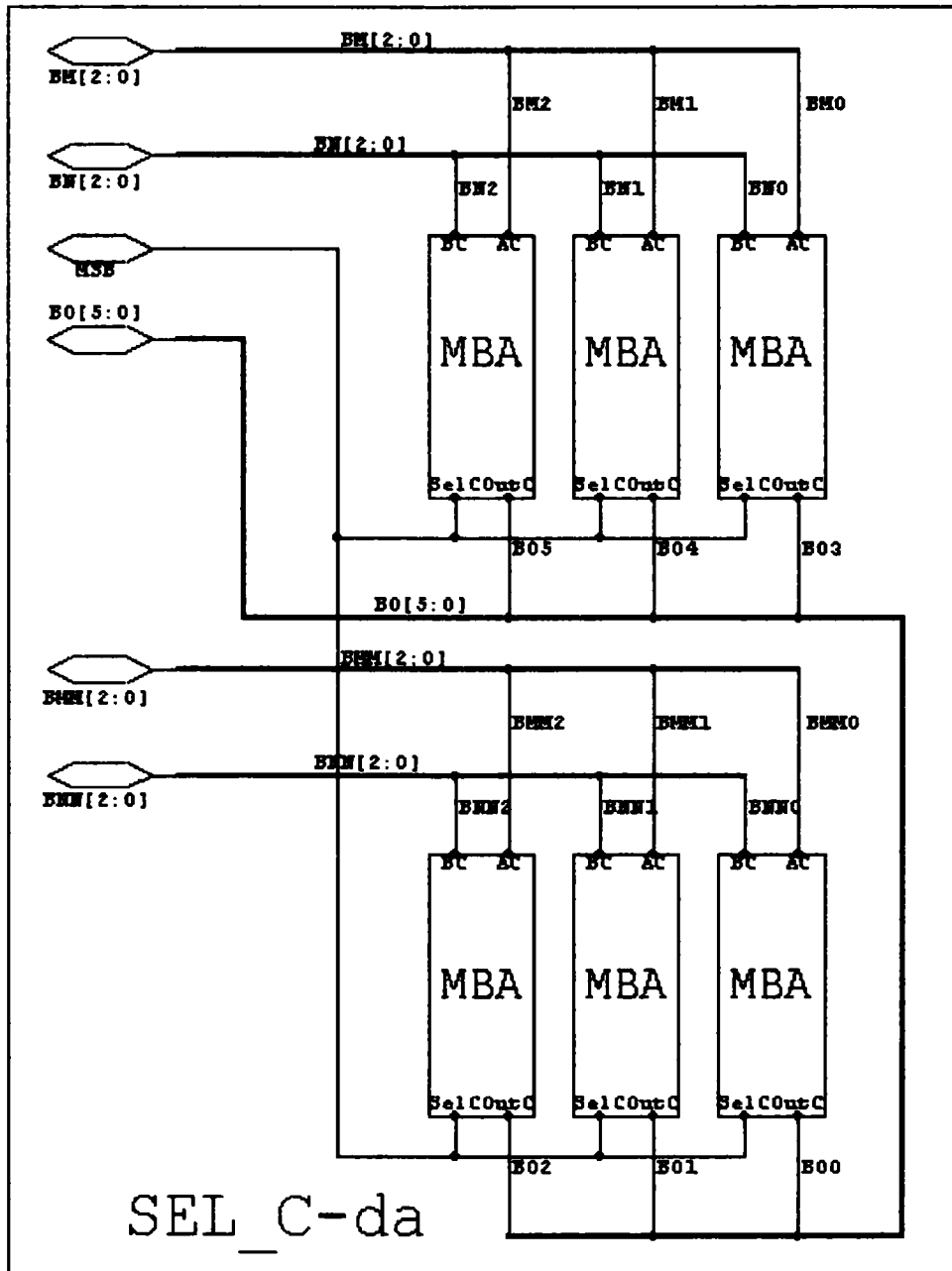
ANEXA 5



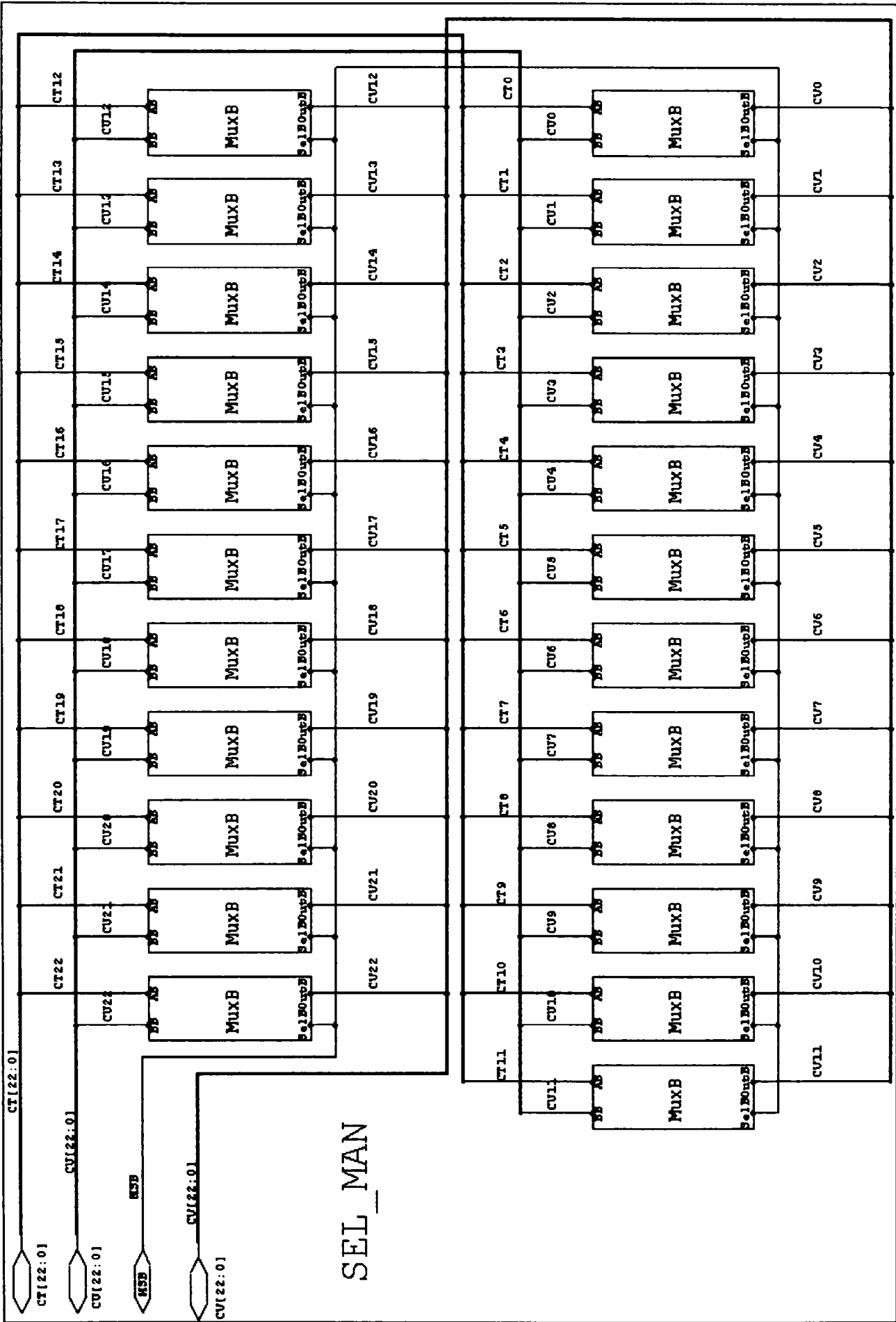
ANEXA 5

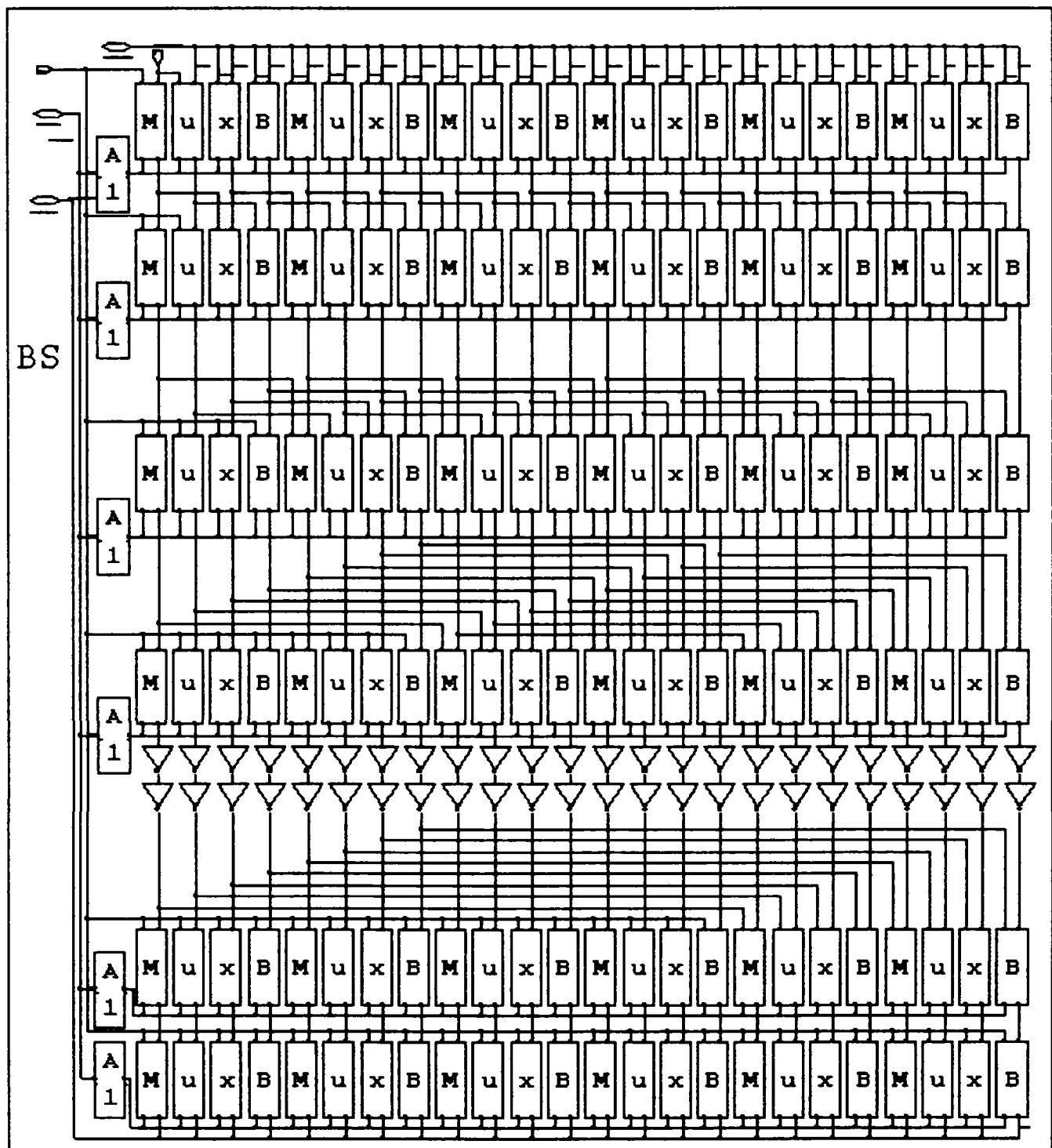


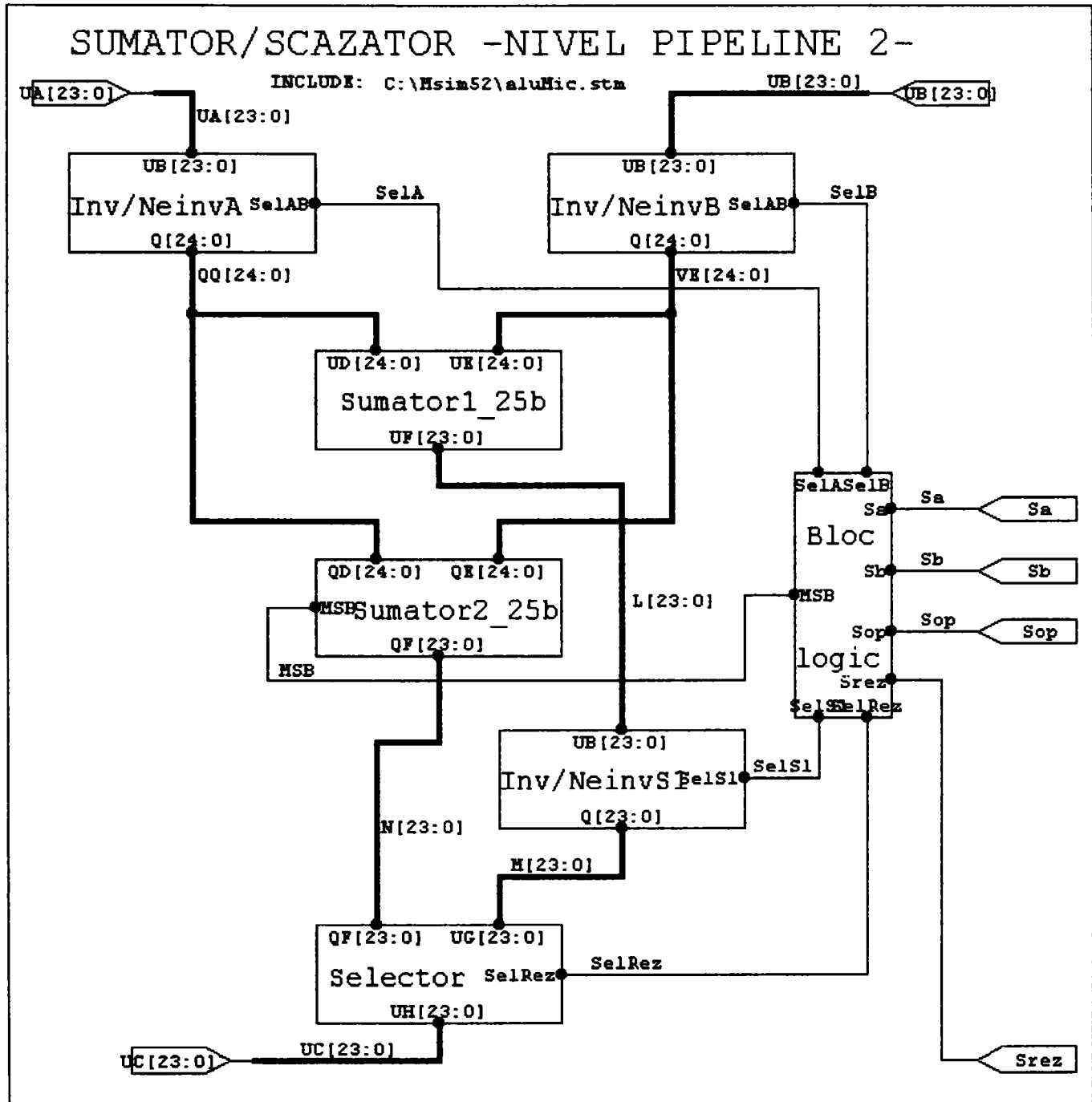
ANEXA 5



ANEXA 5

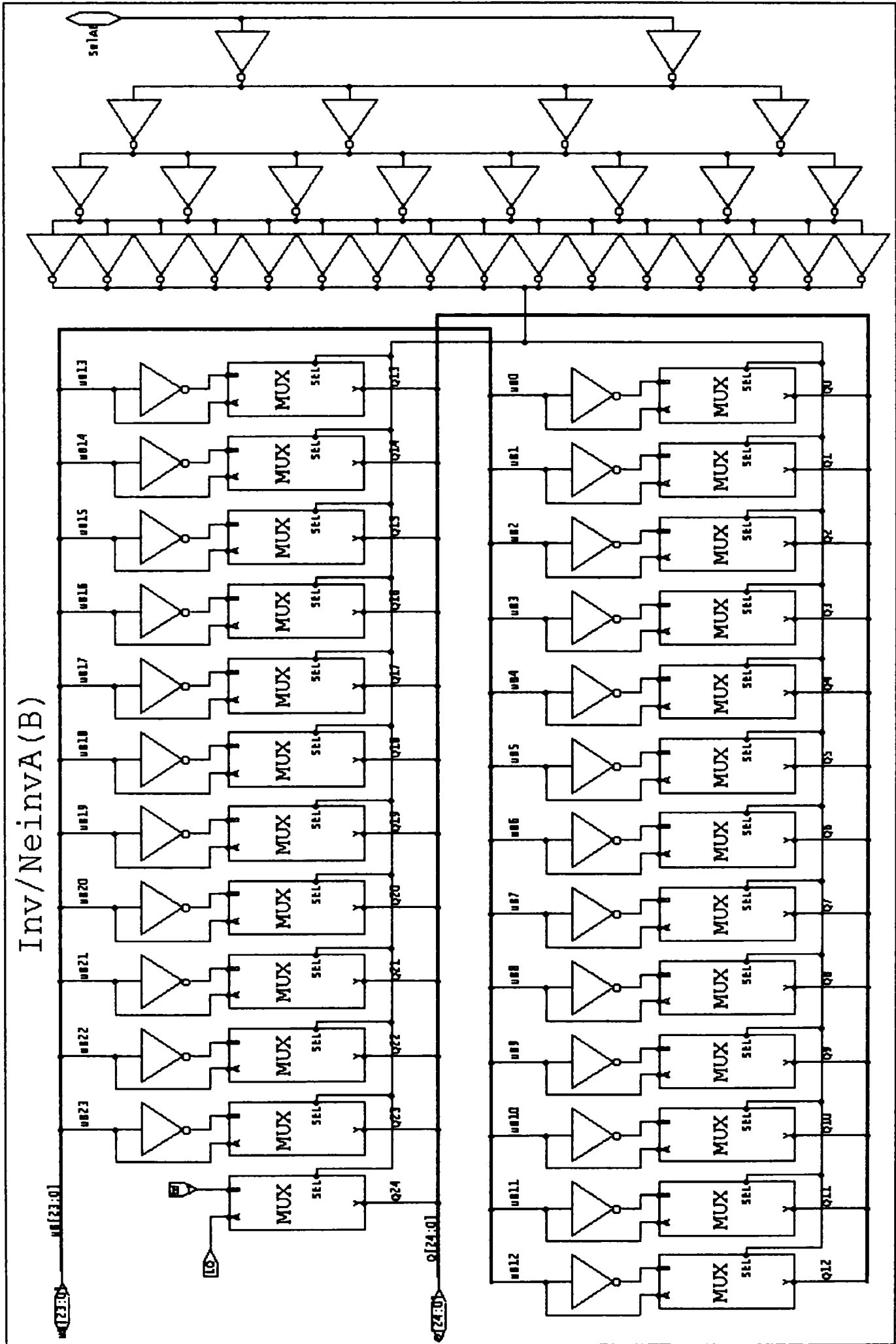


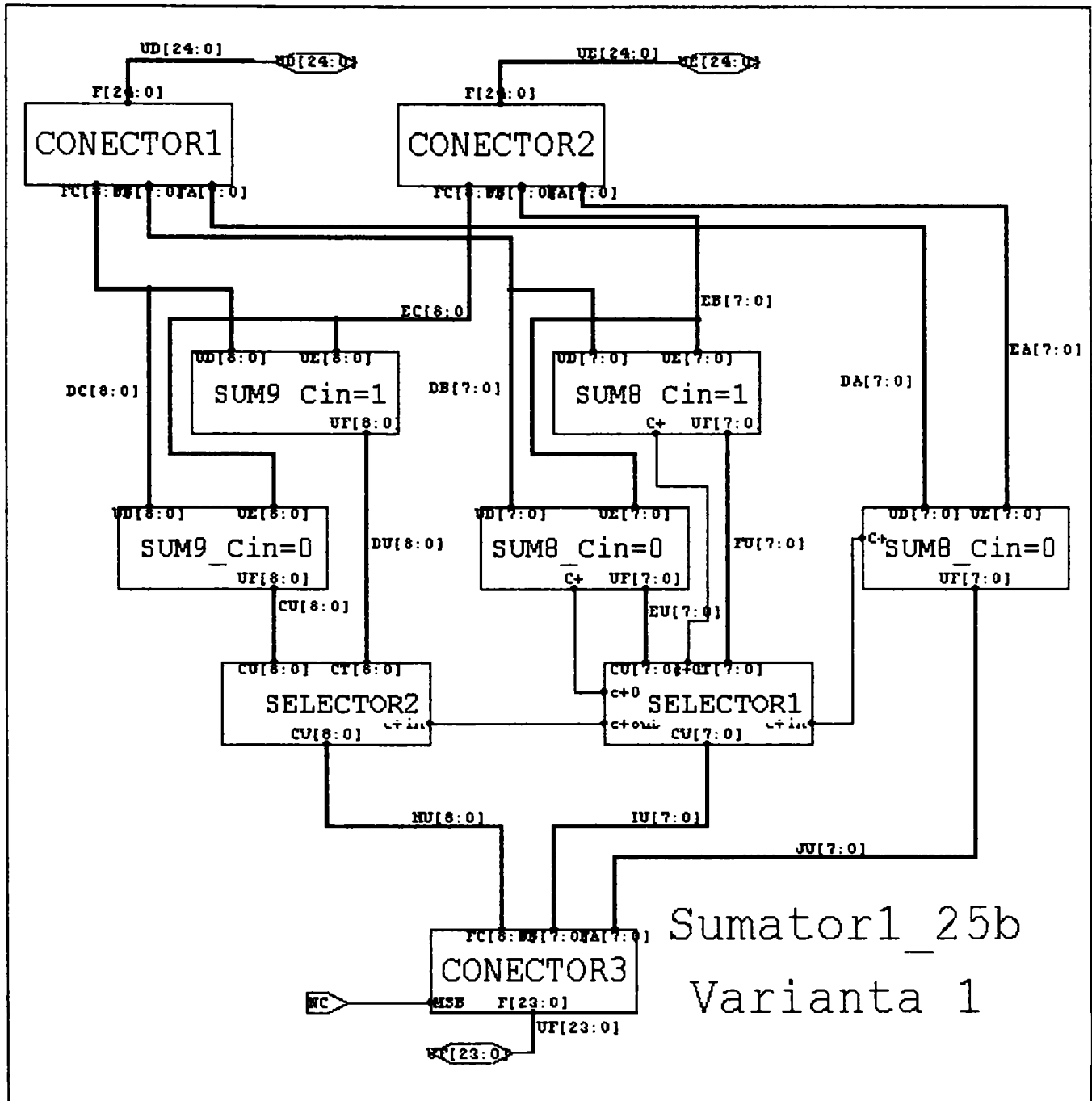


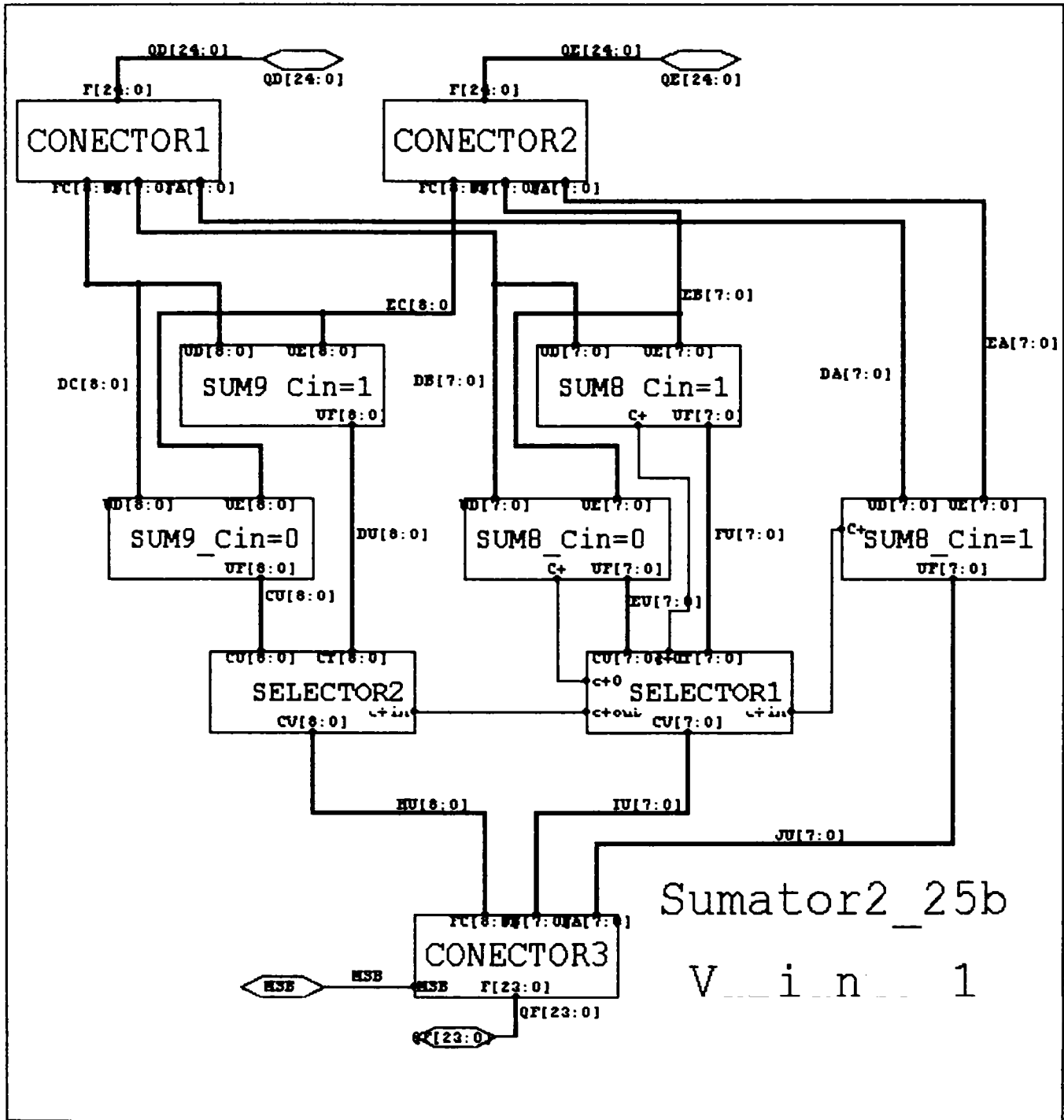


ANEXA 6

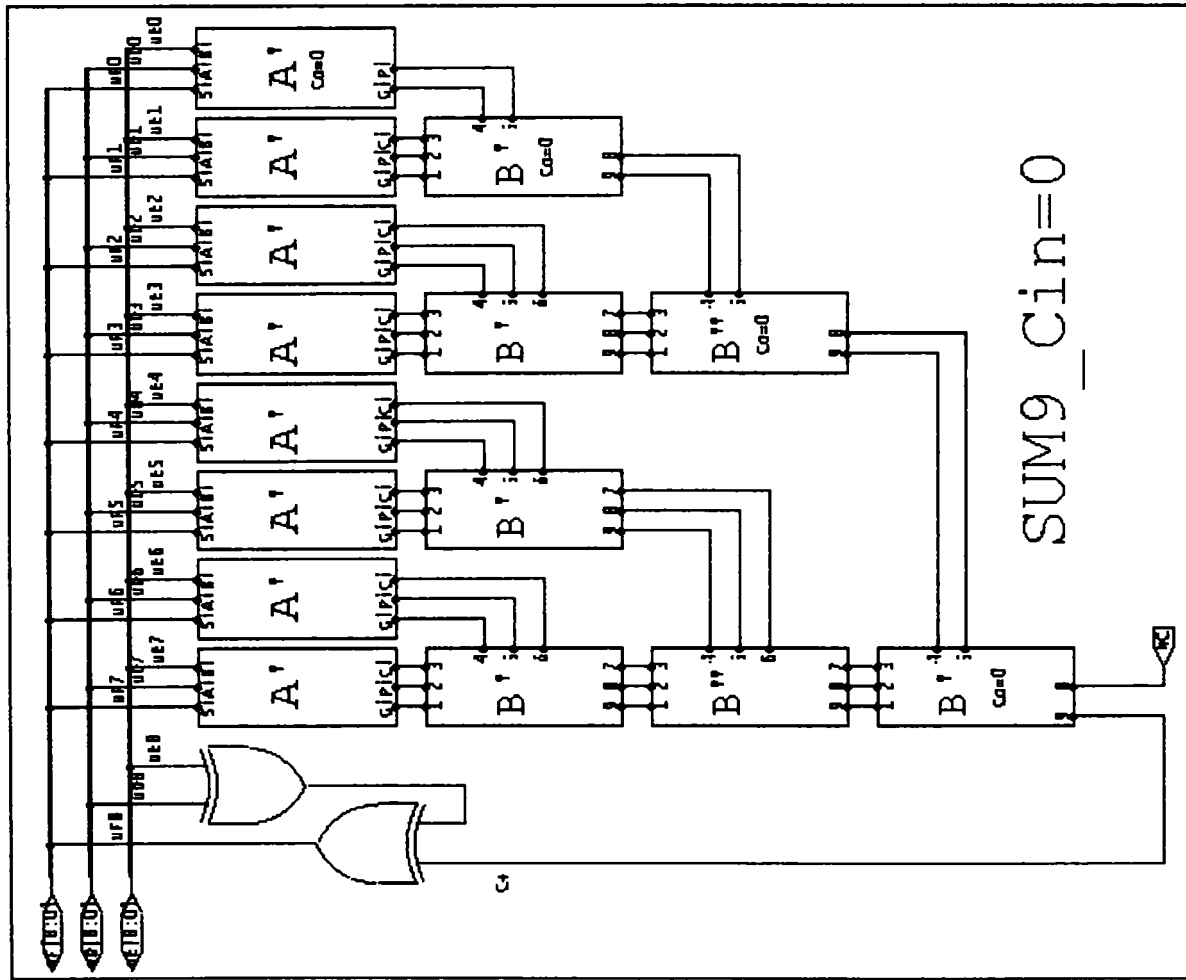
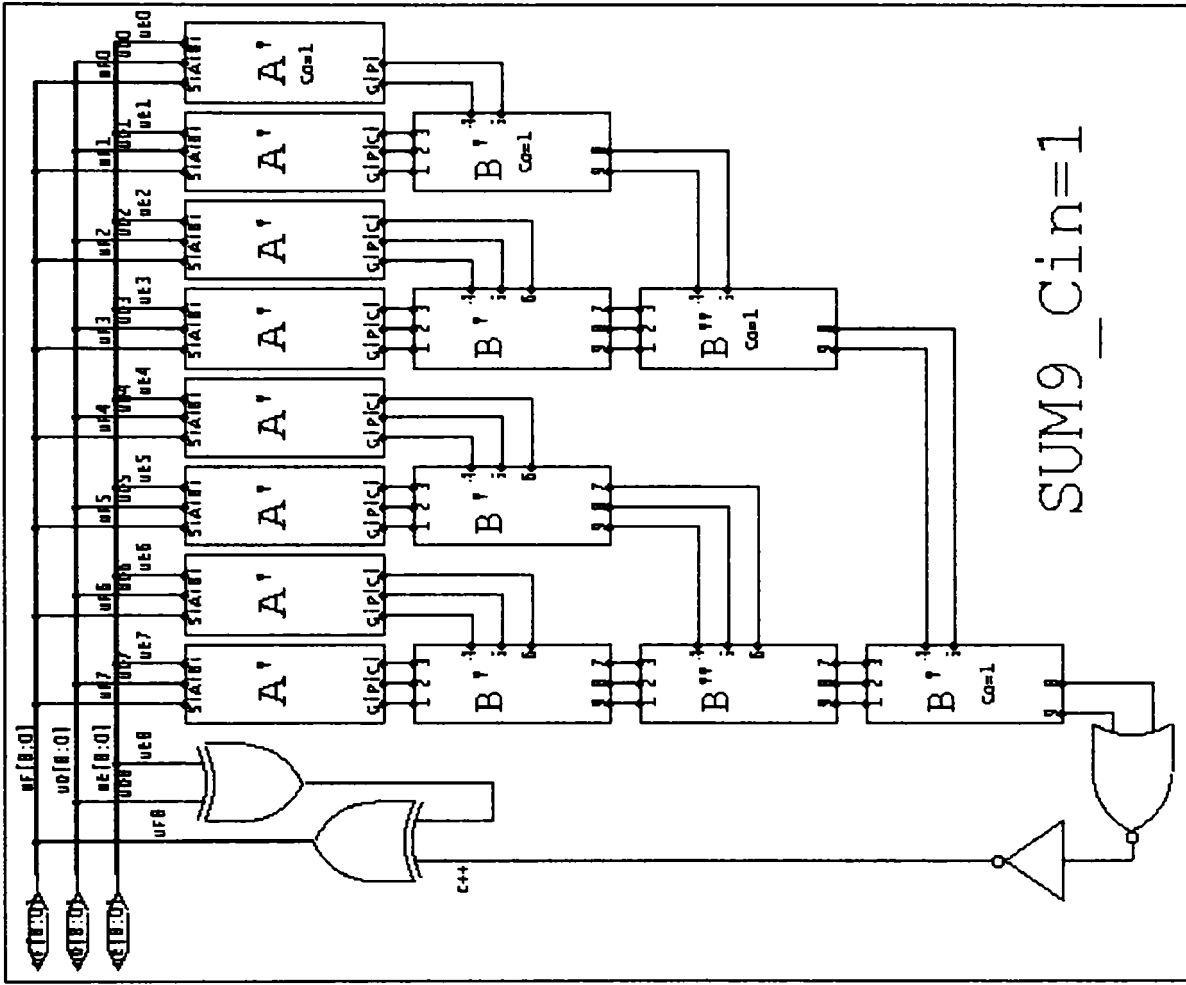
Inv/NeinVA(B)



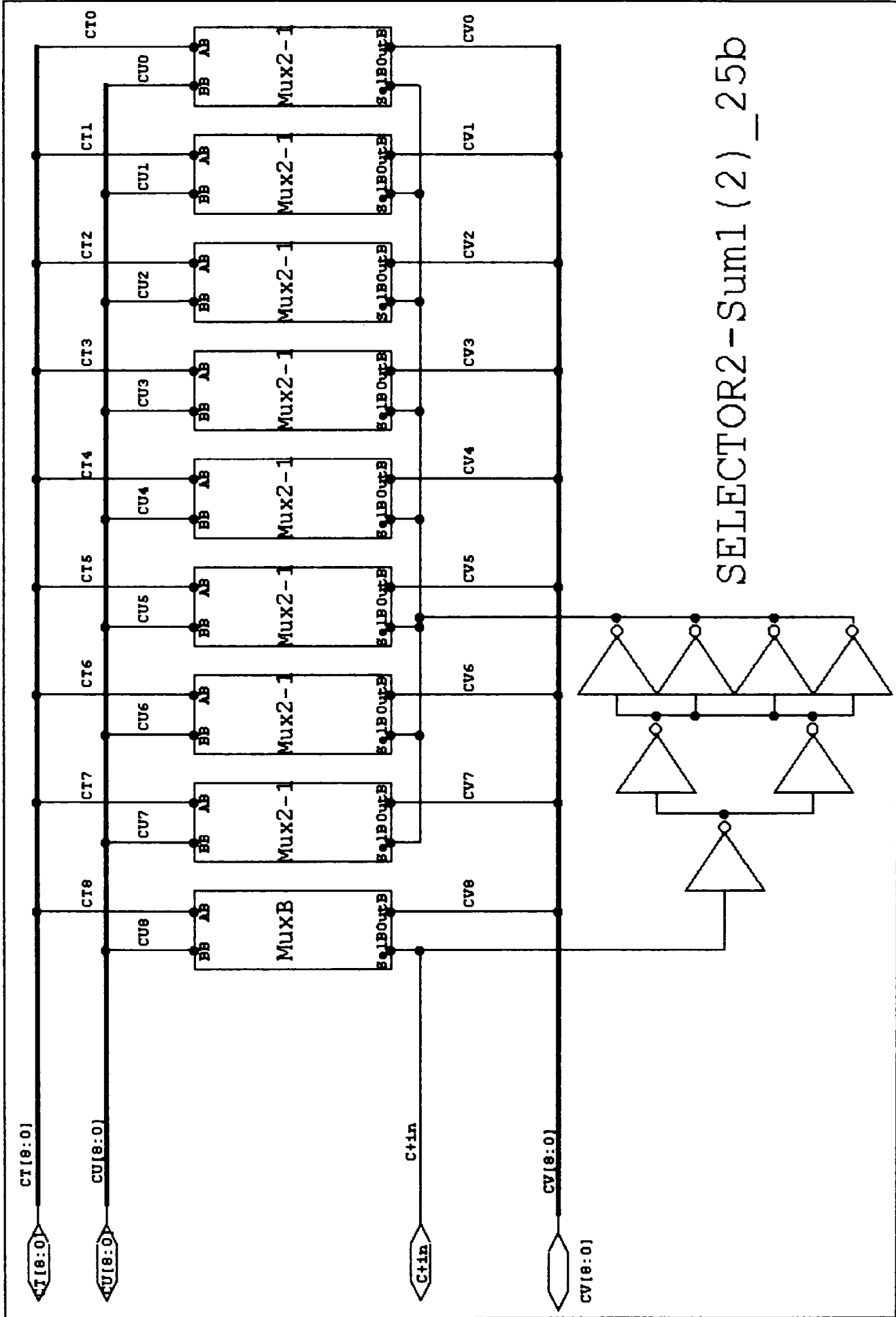




ANEXA 6

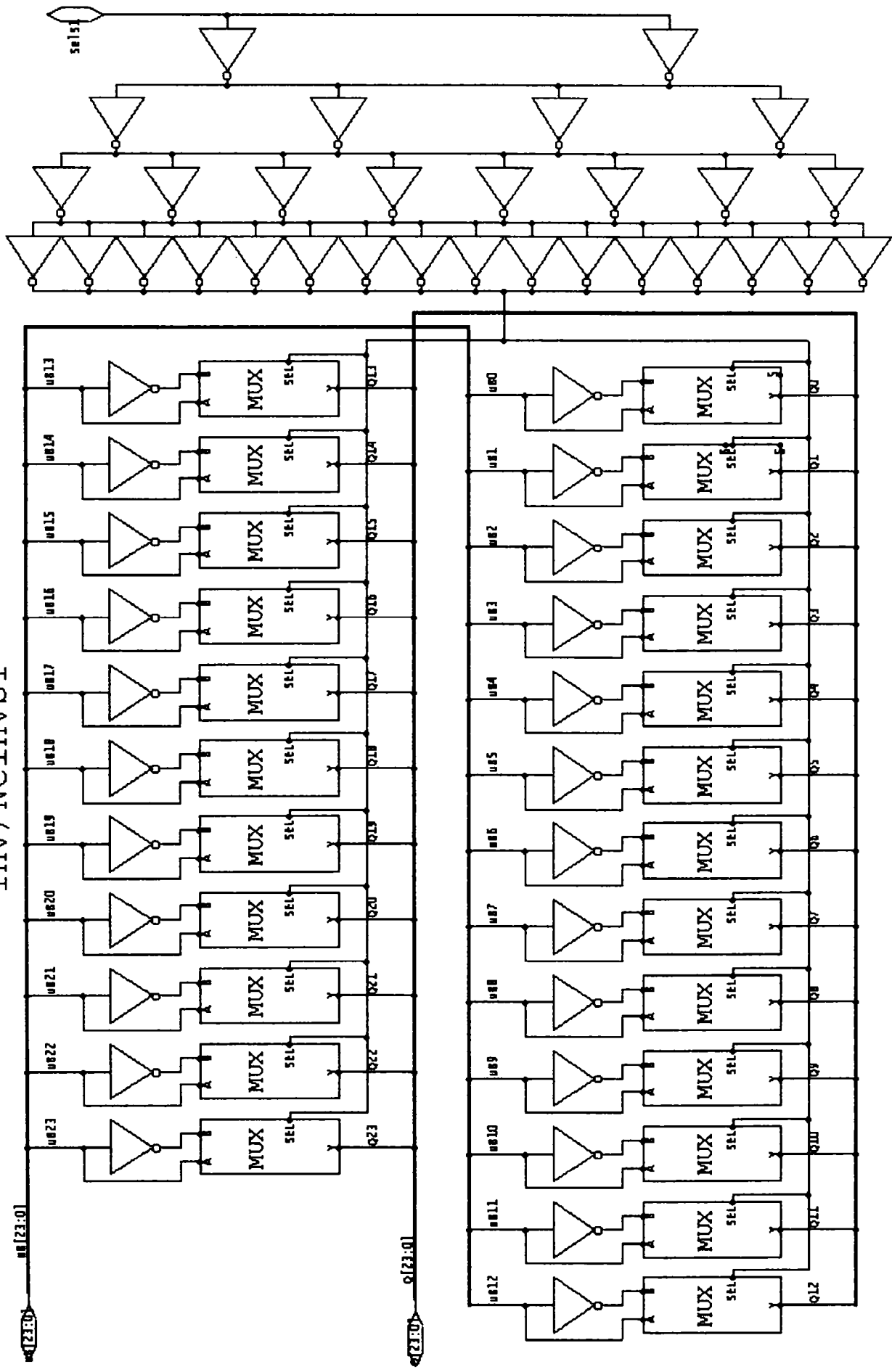


ANEXA 6



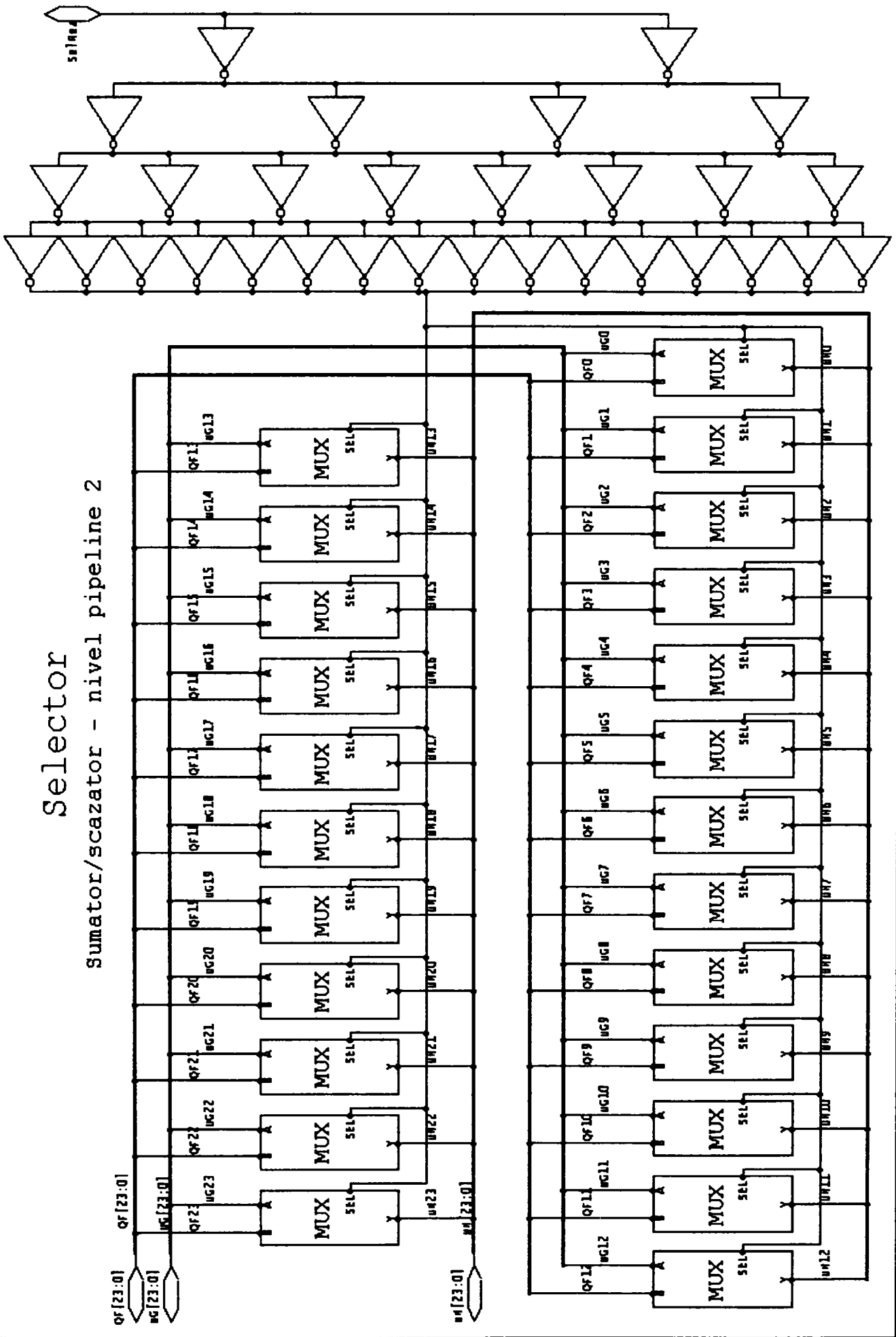
ANEXA 6

Inv/NeInvS1

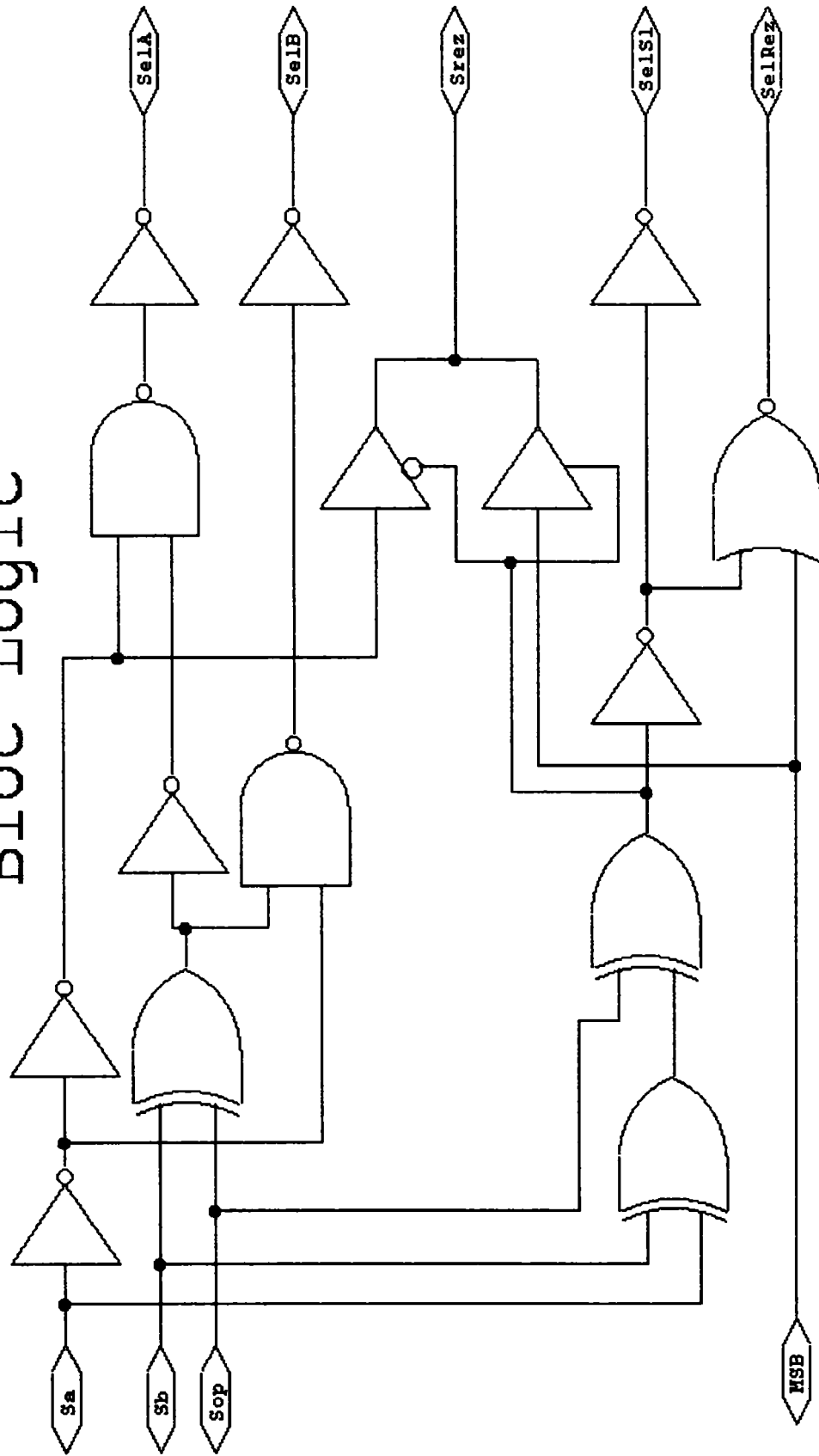


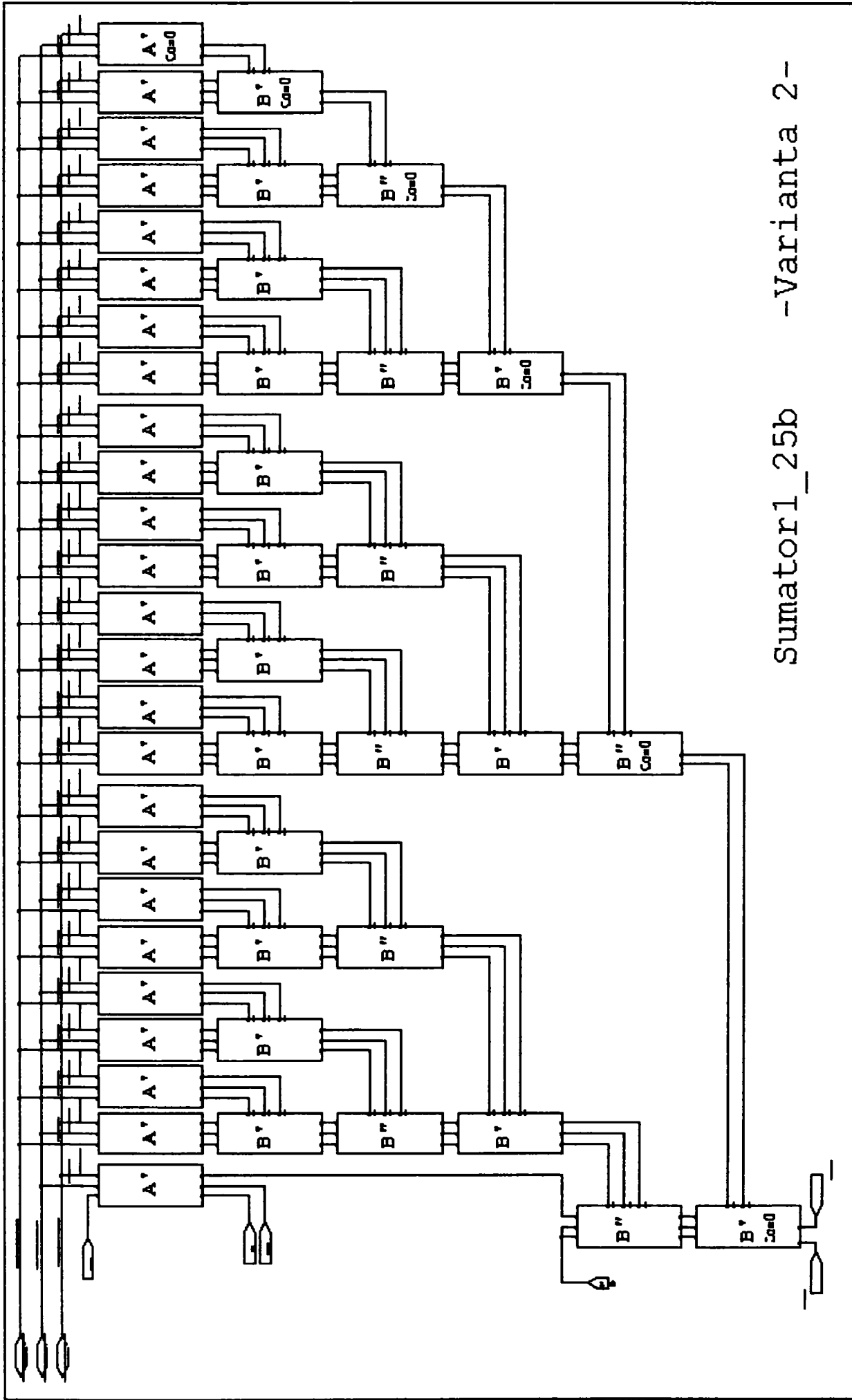
A6.7

Selector Sumator/scazator - nivel pipeline 2

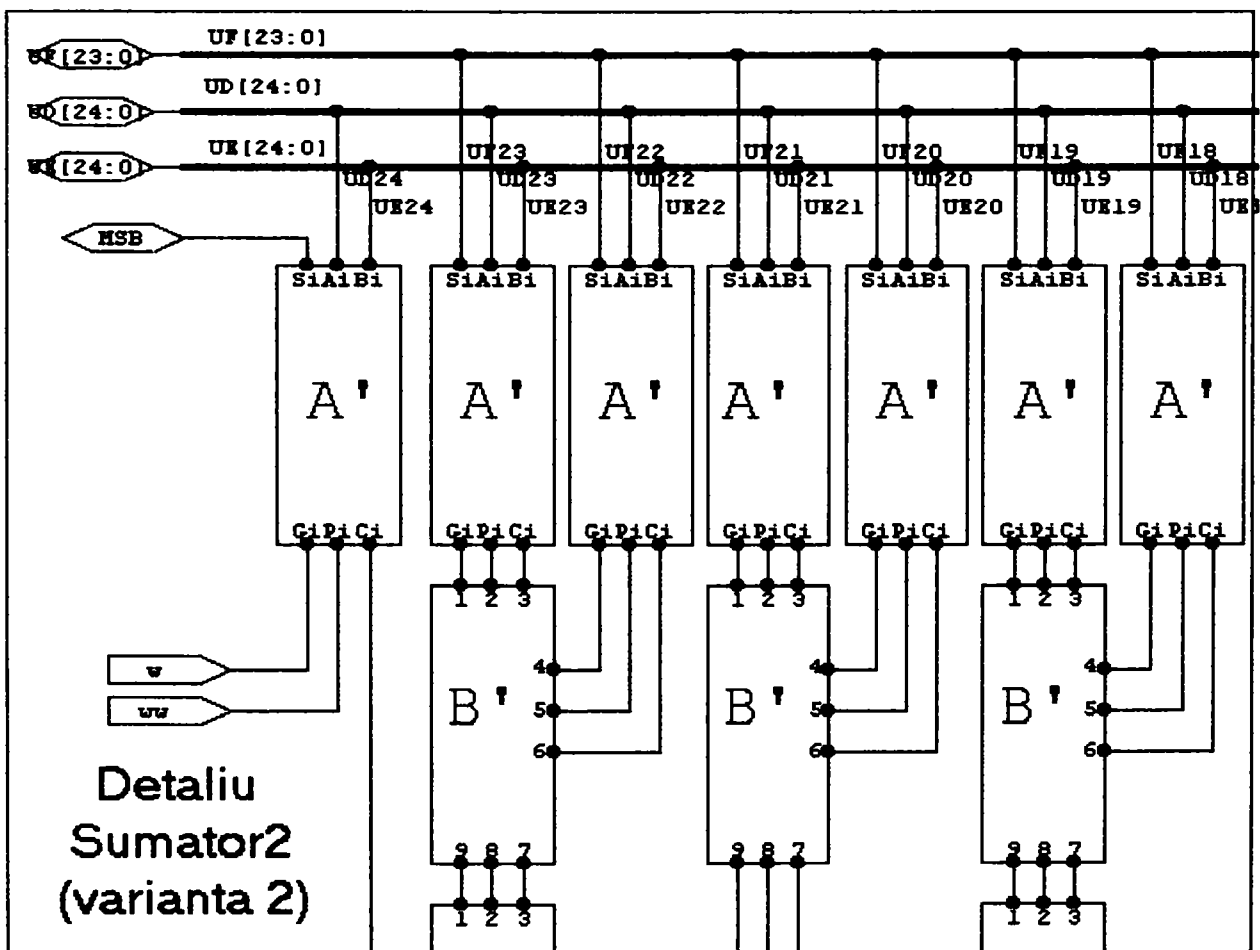
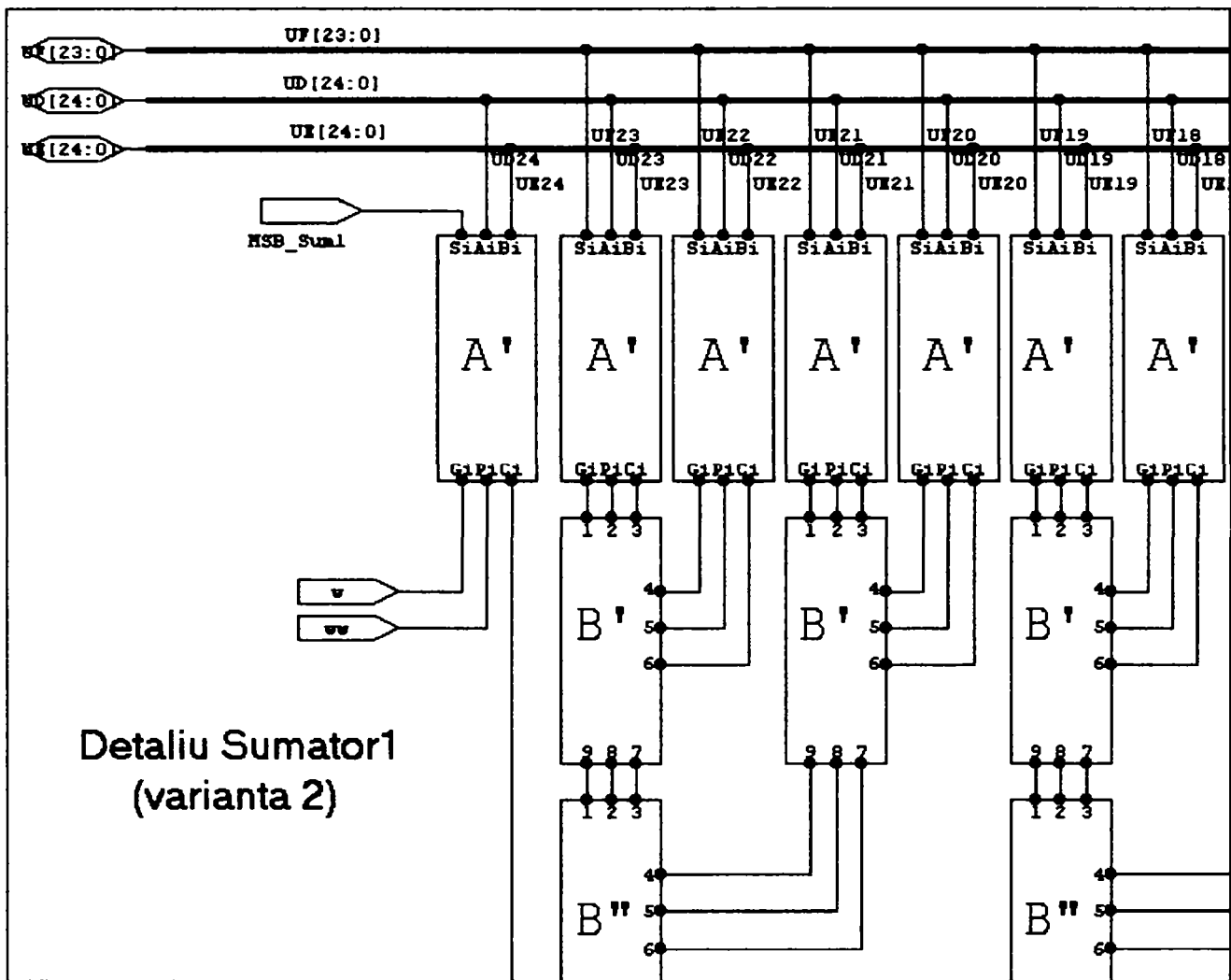


Bloc Logic

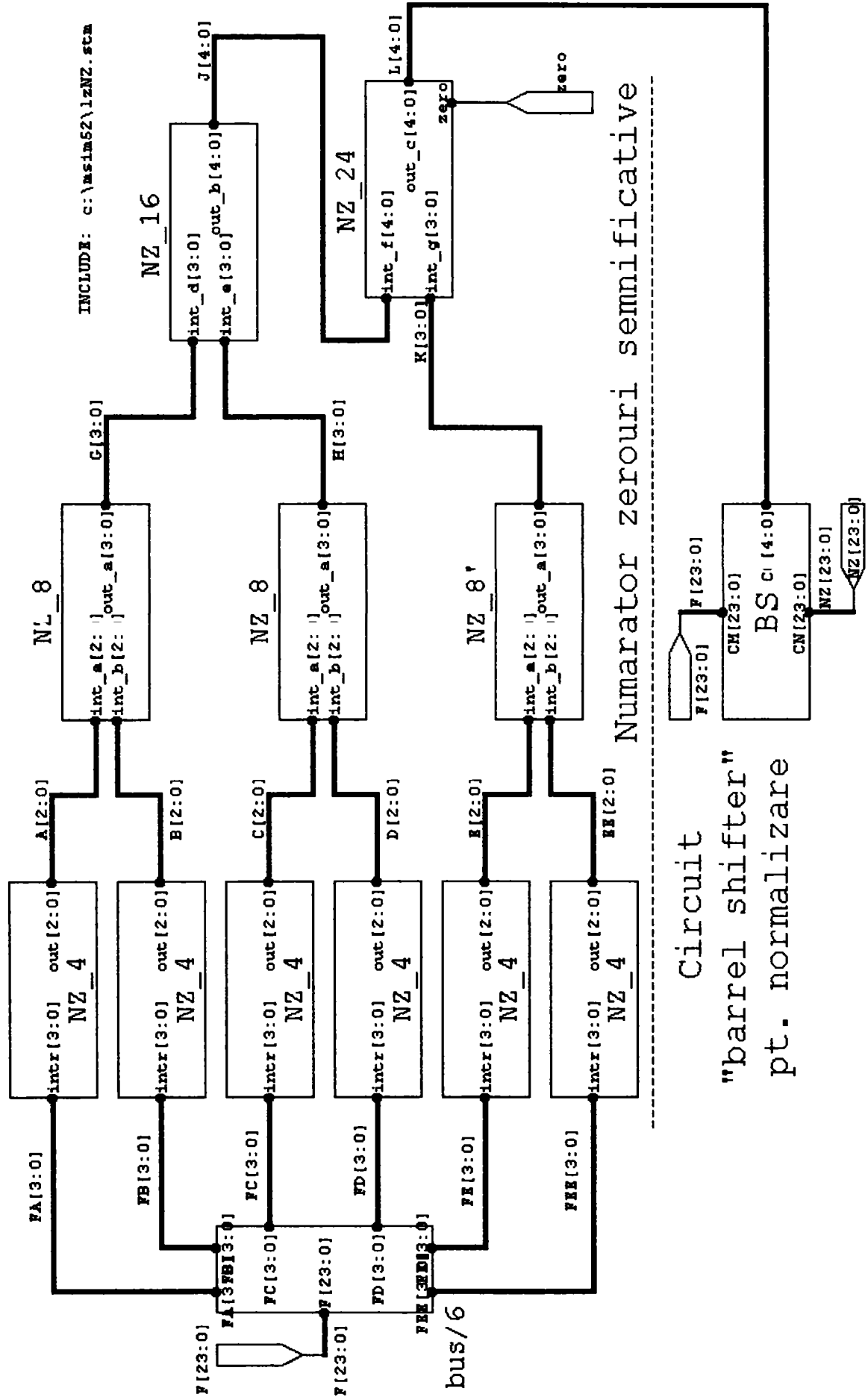




Sumator1_25b -Varianta 2-

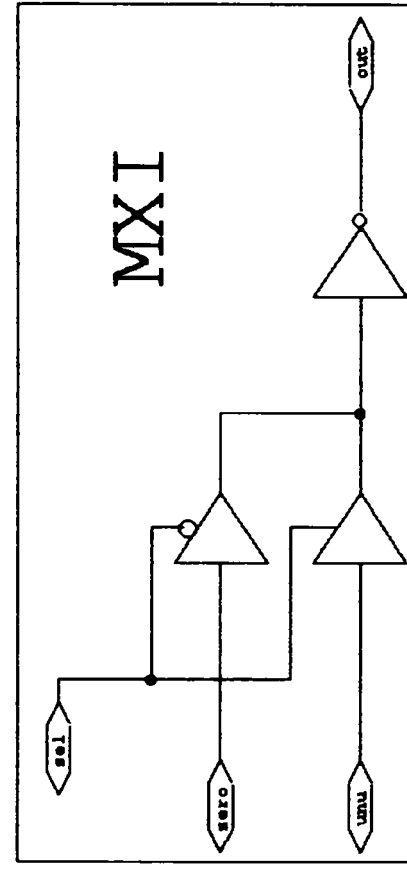
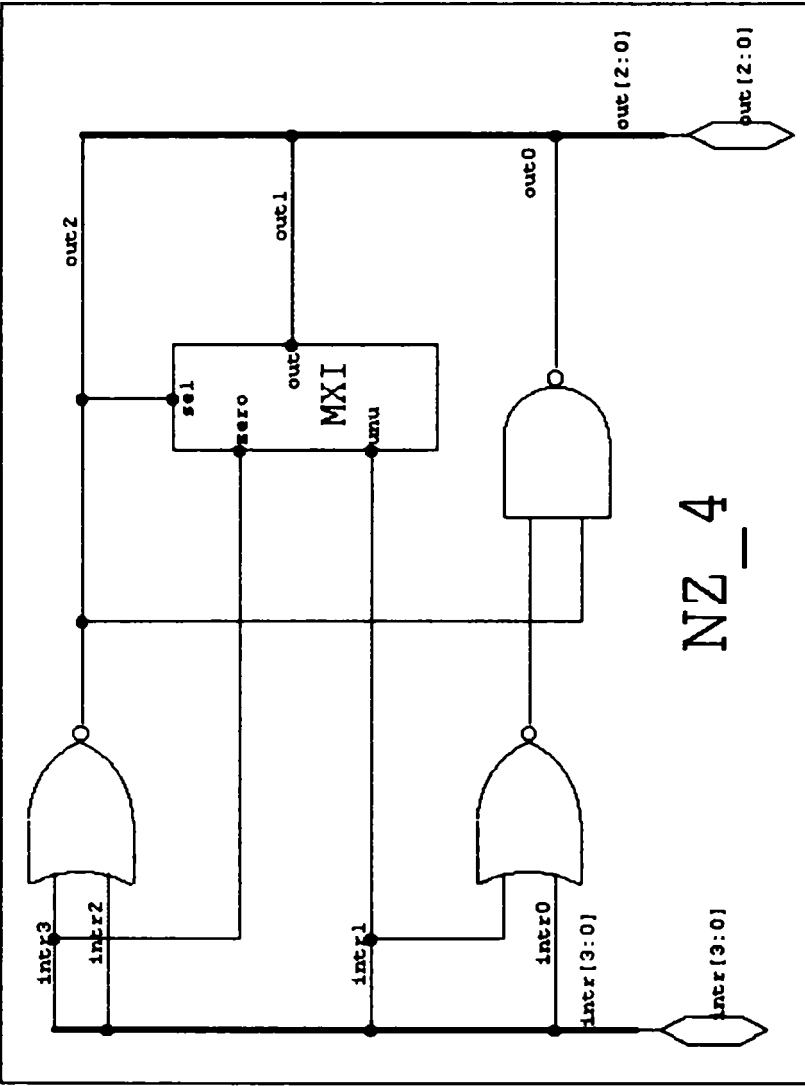
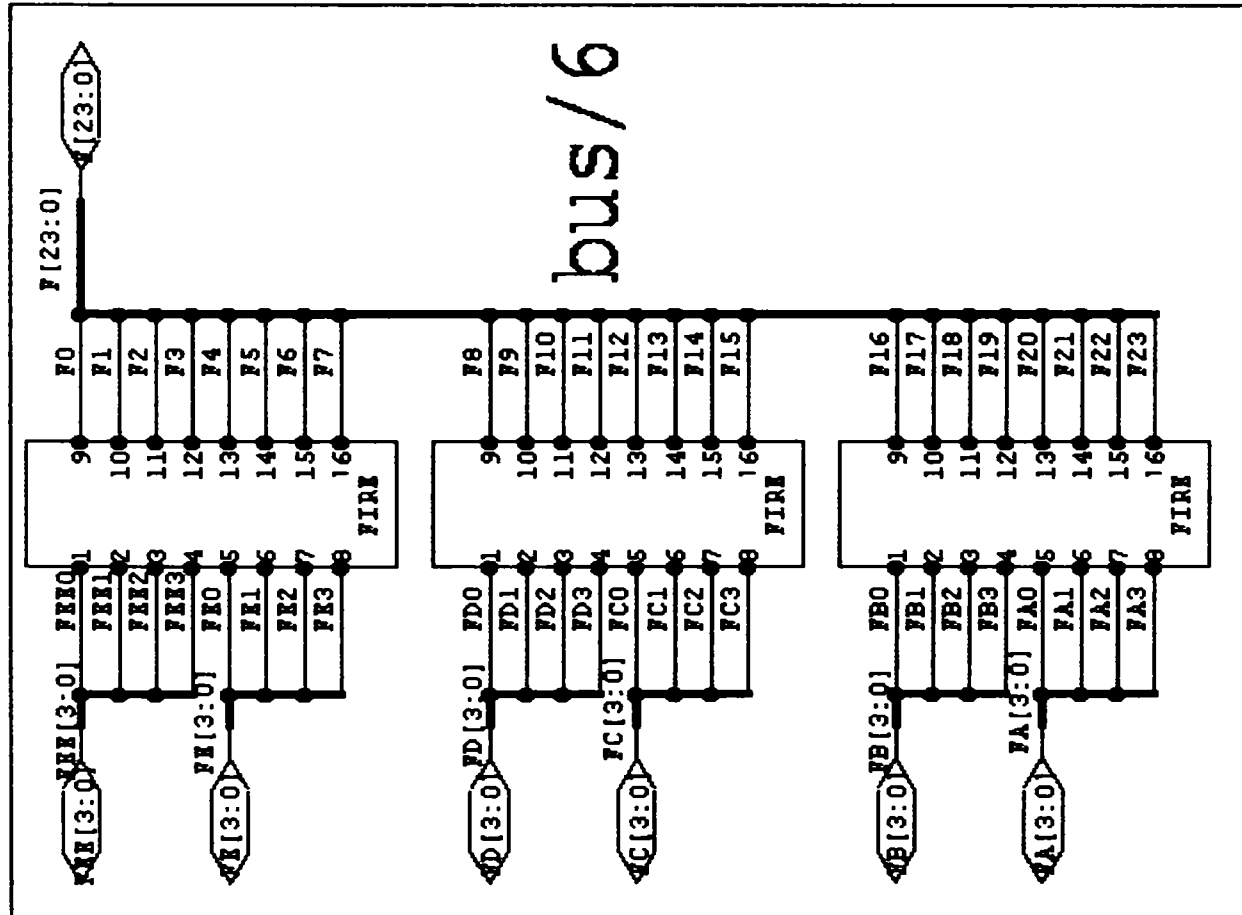


CIRCUIT DE NORMALIZARE

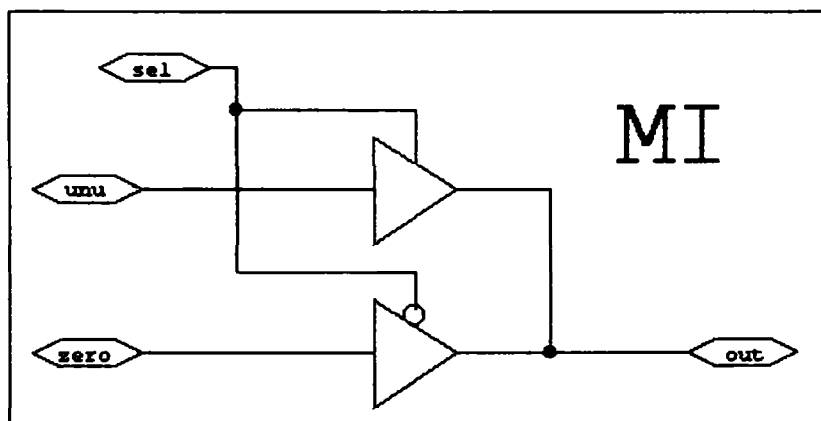
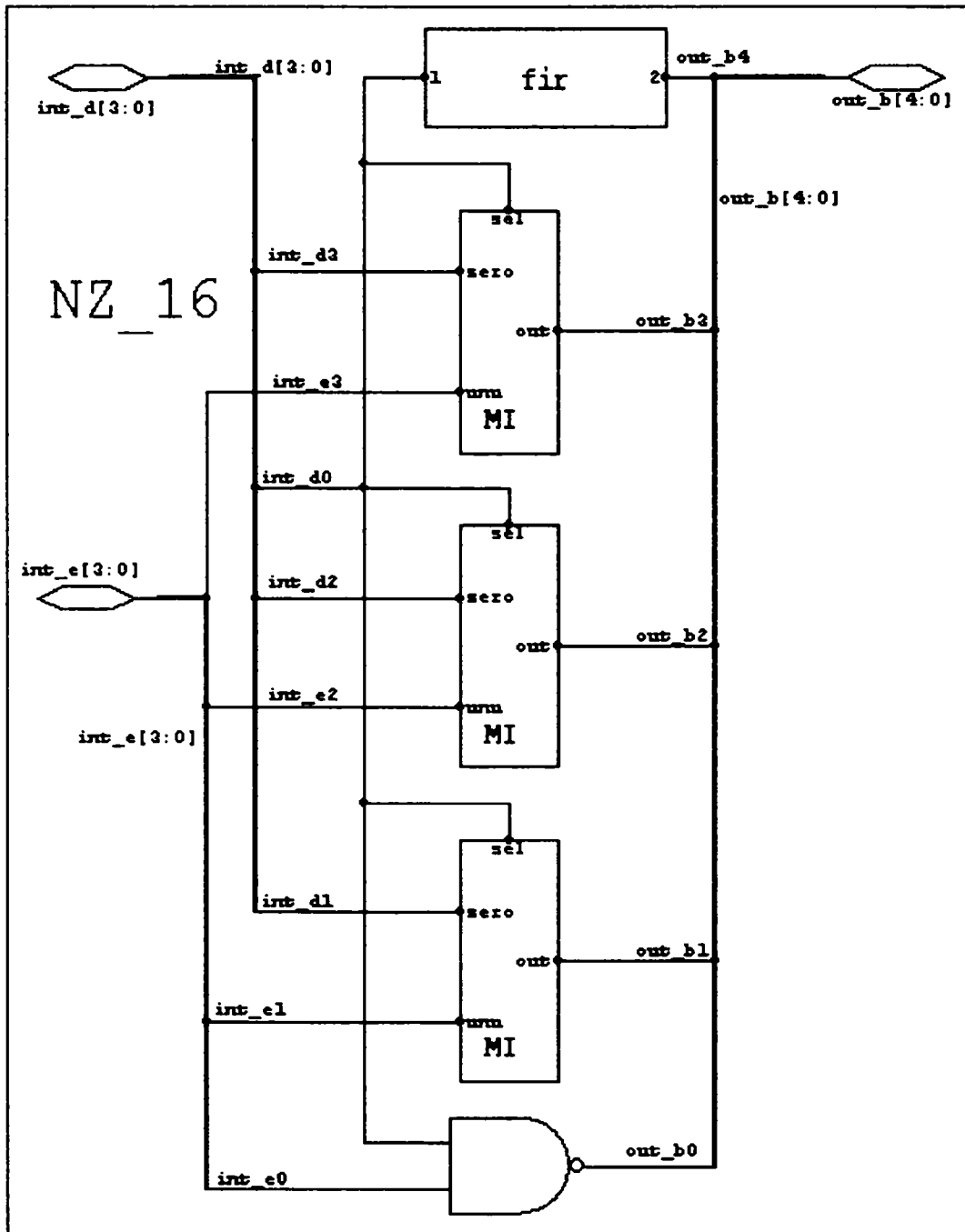


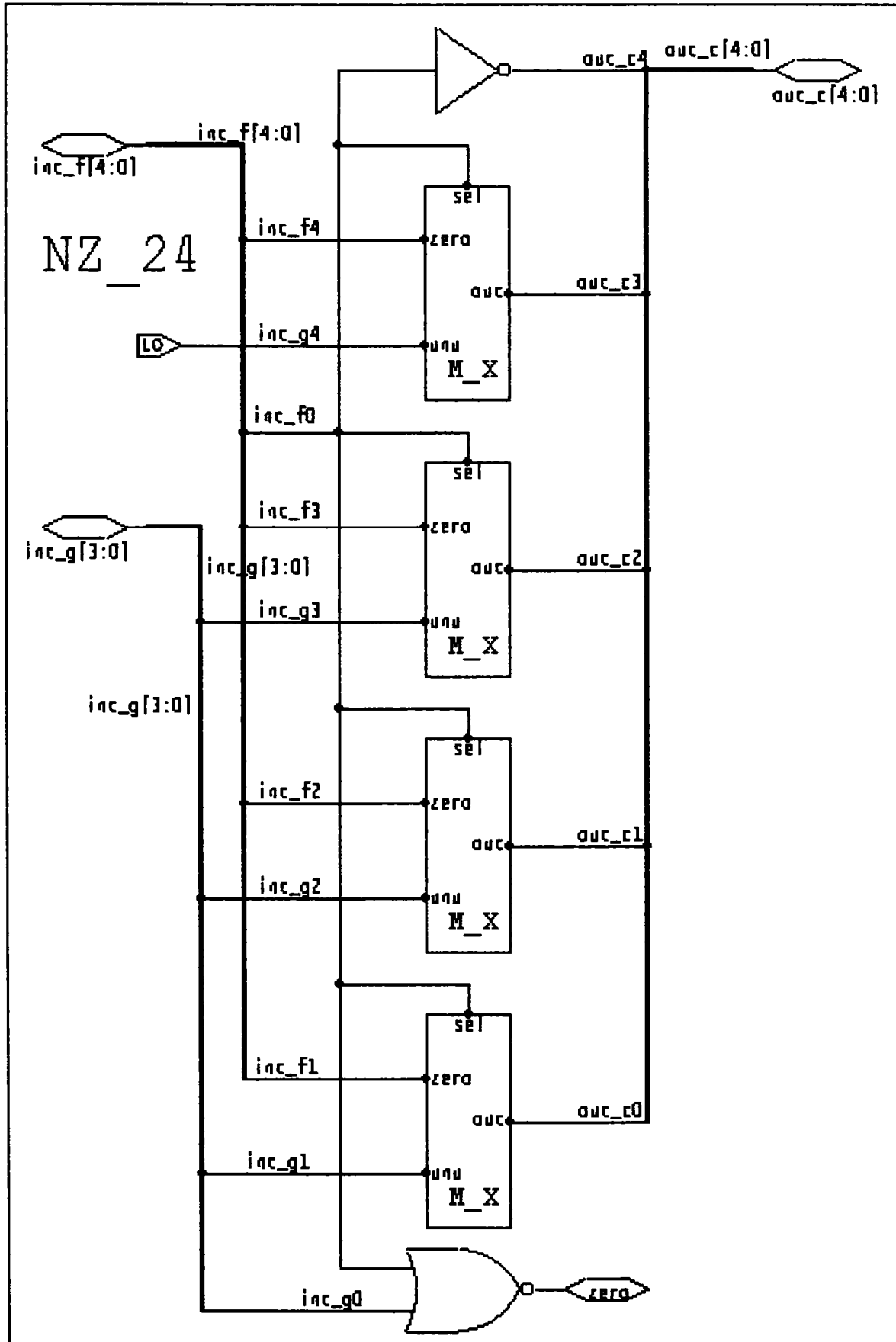
Numarator zerouri semnificative

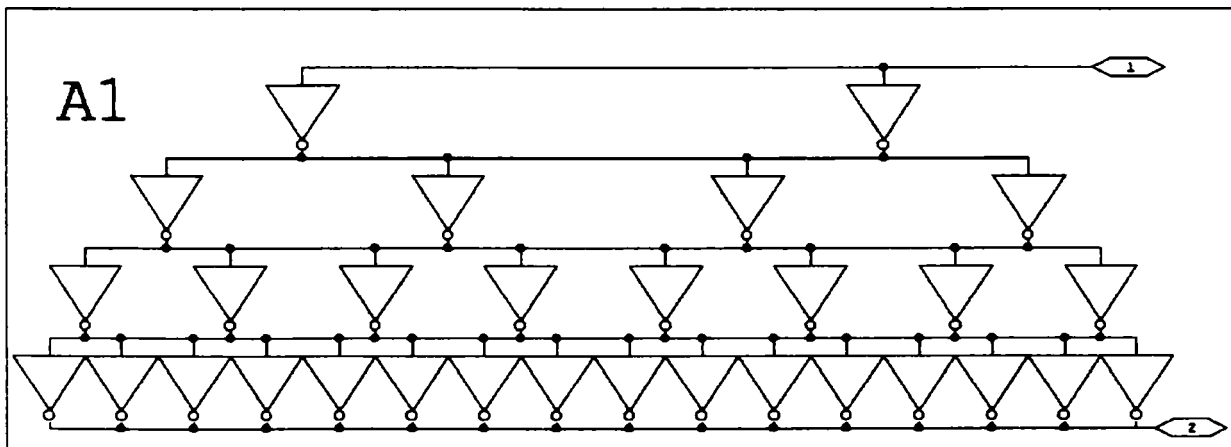
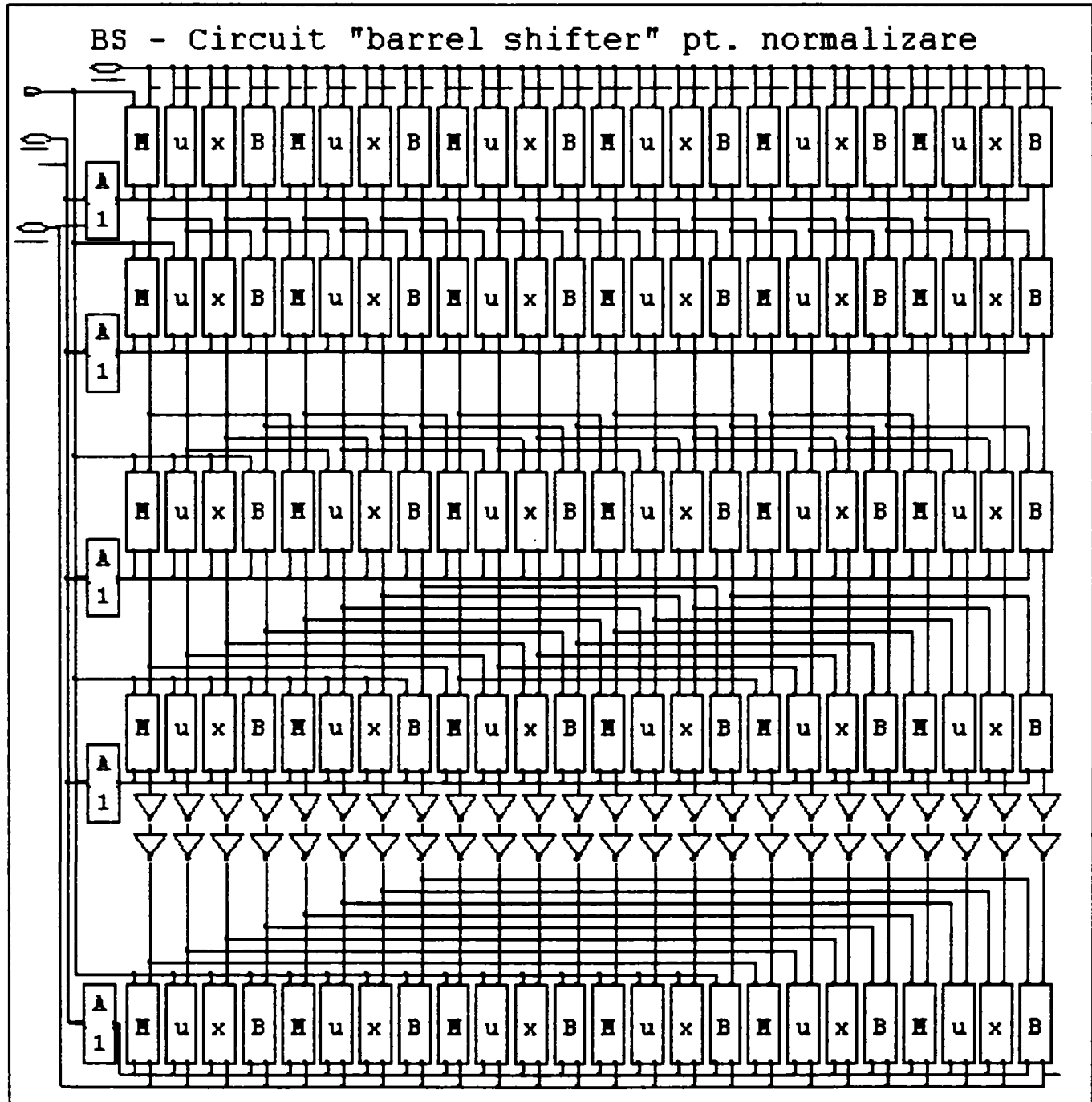
Circuit
"barrel shifter"
pt. normalizare

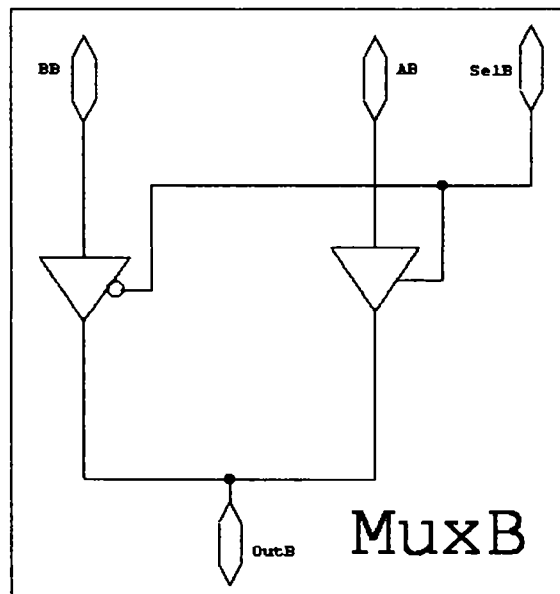
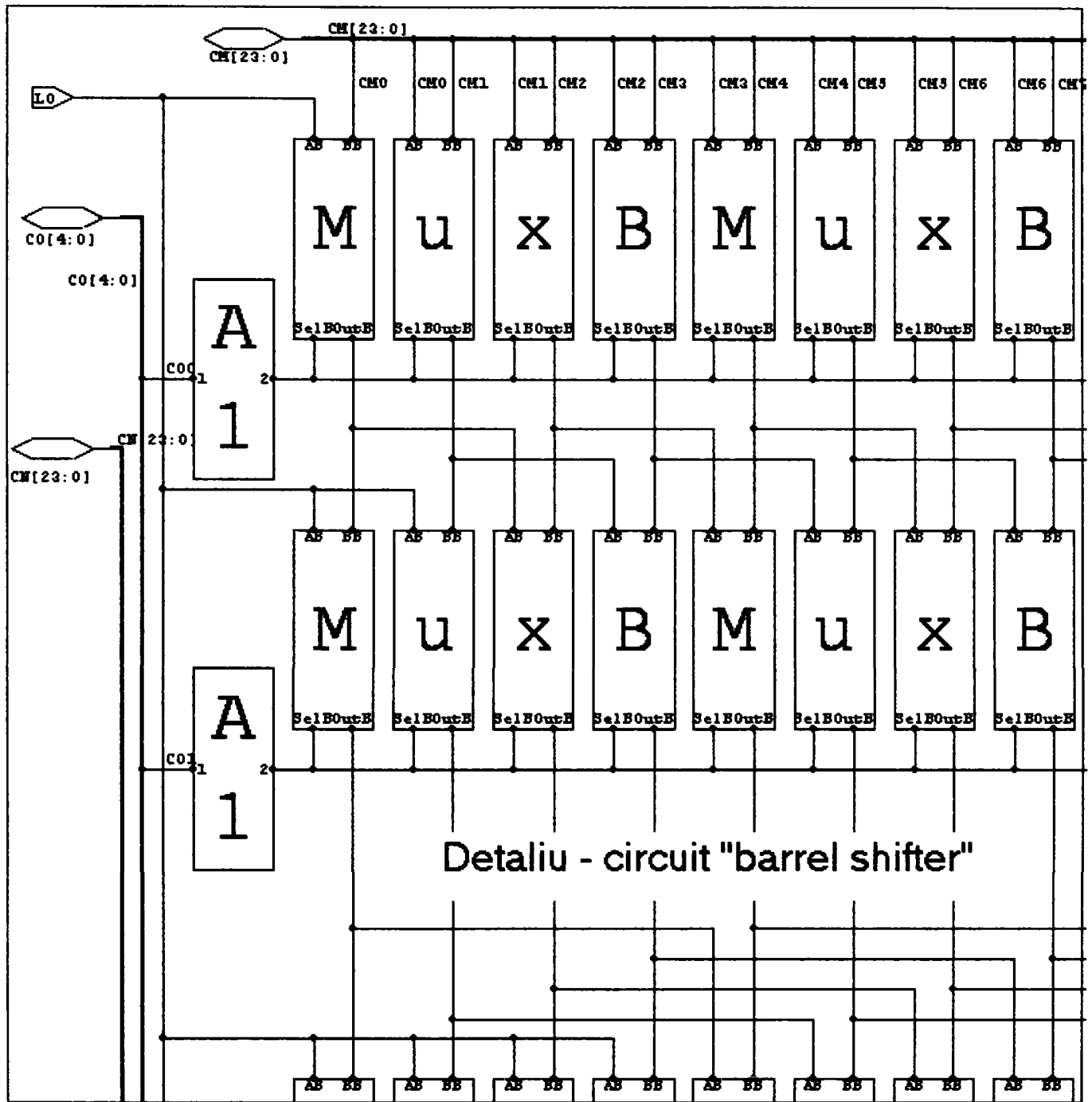


ANEXA 7





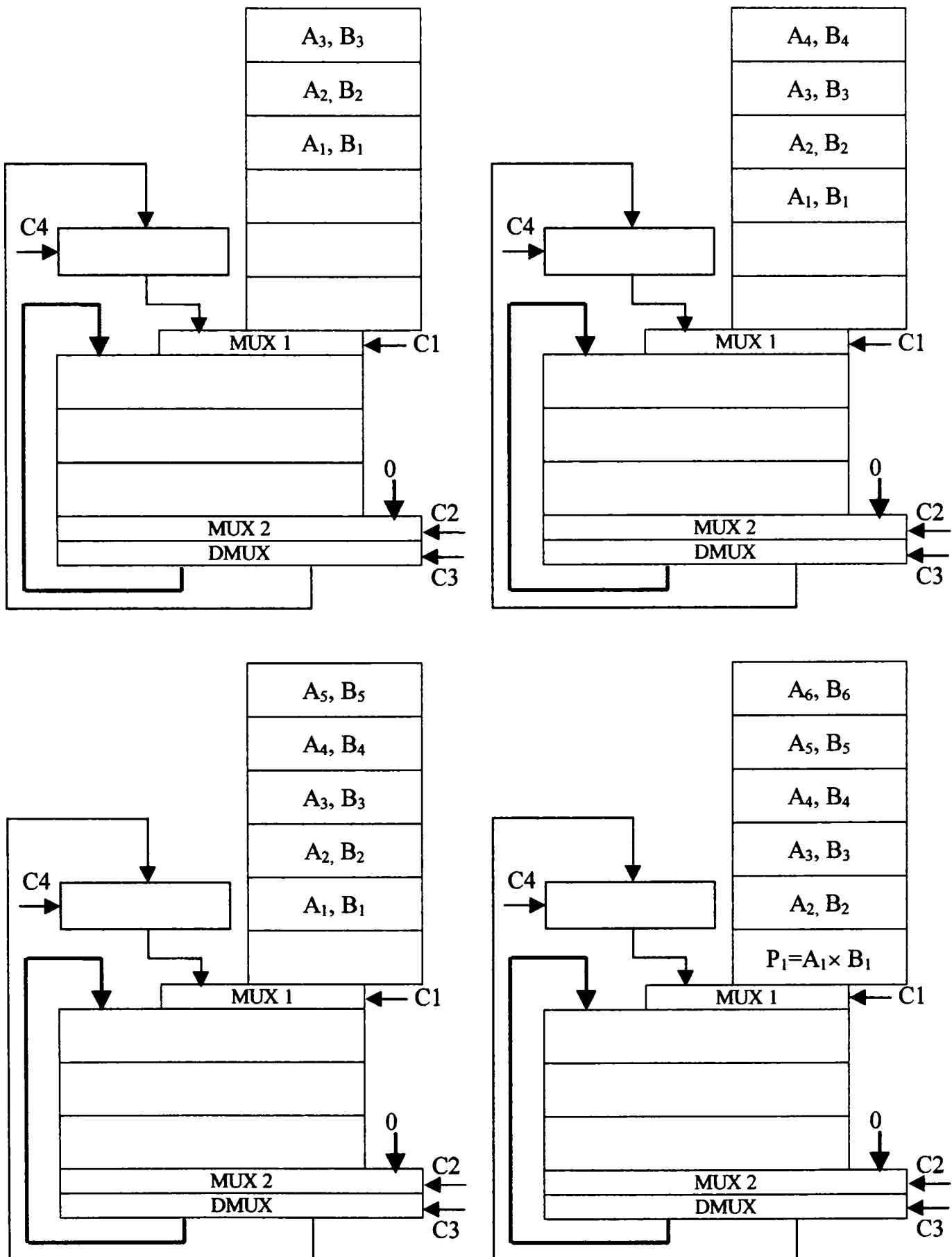




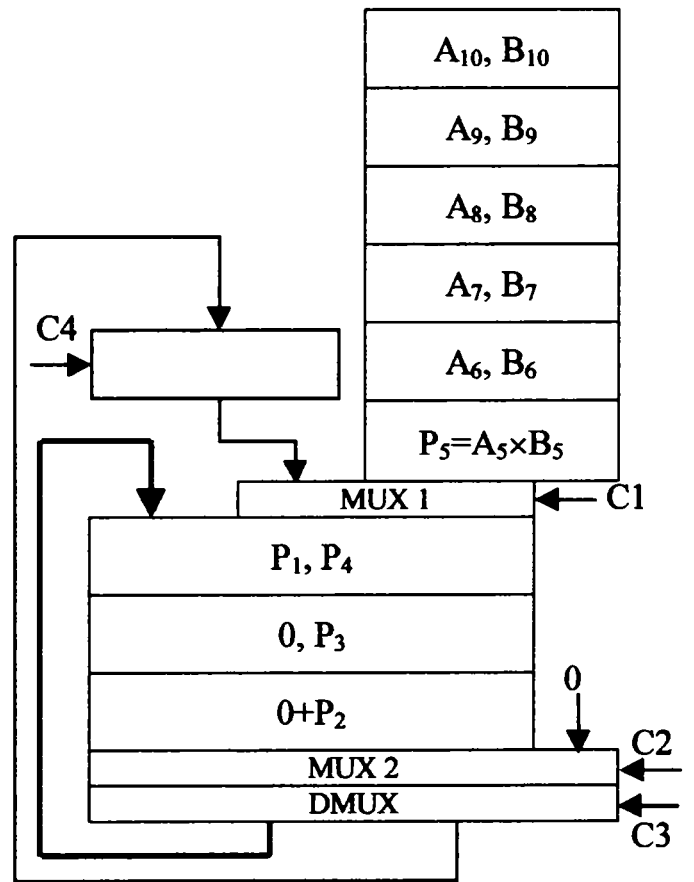
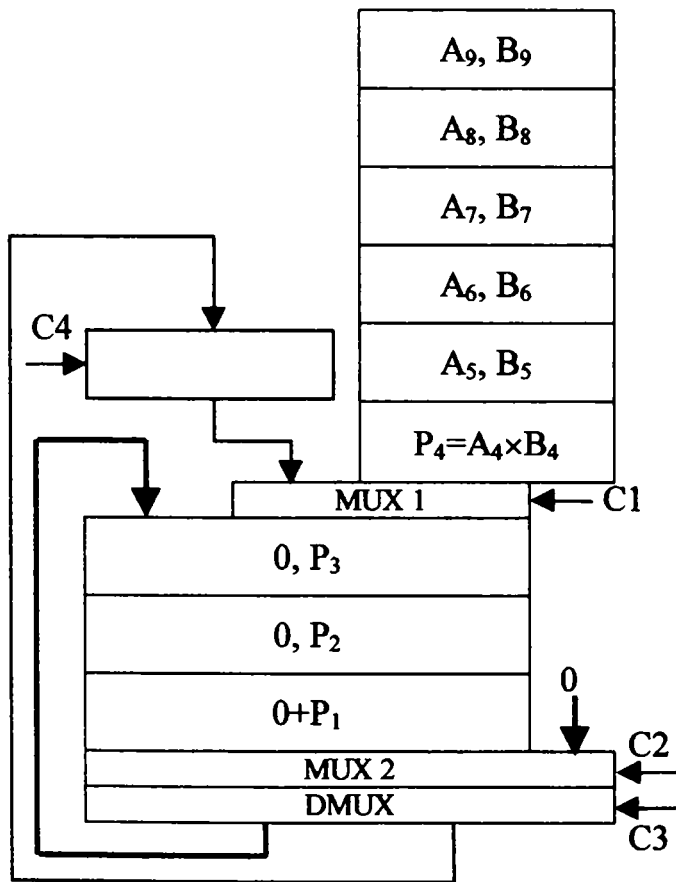
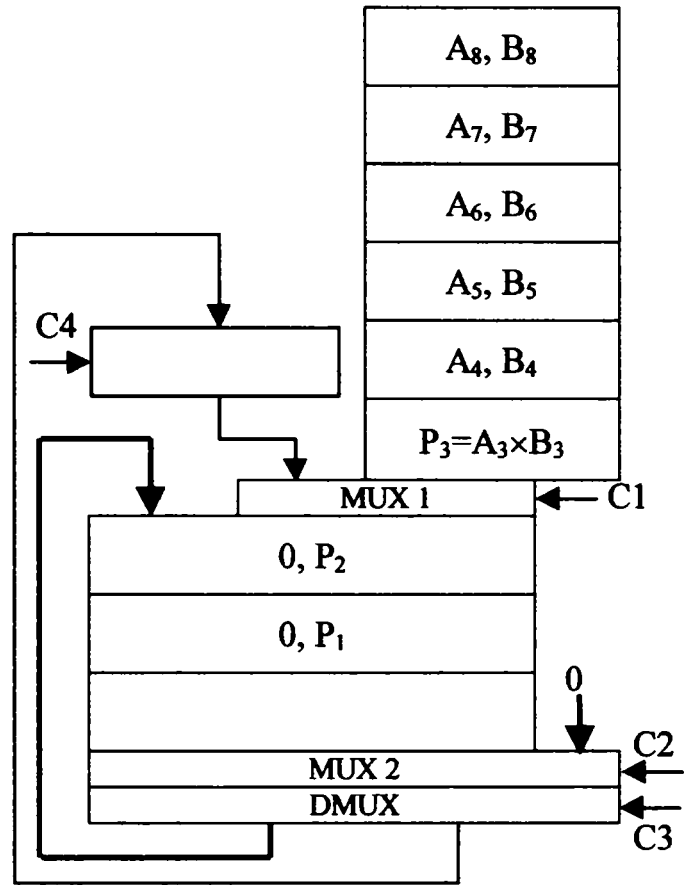
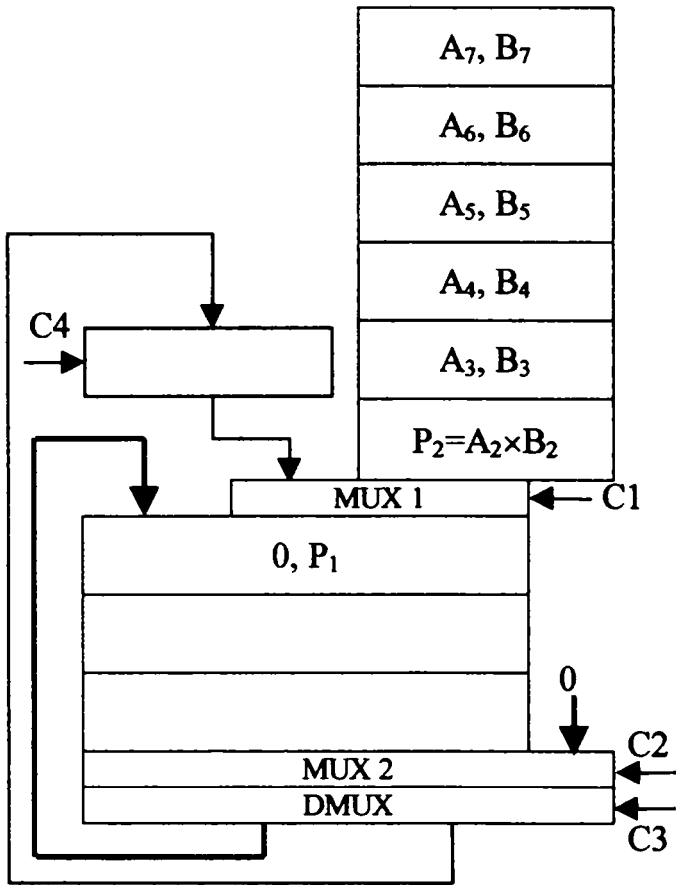
ANEXA 8

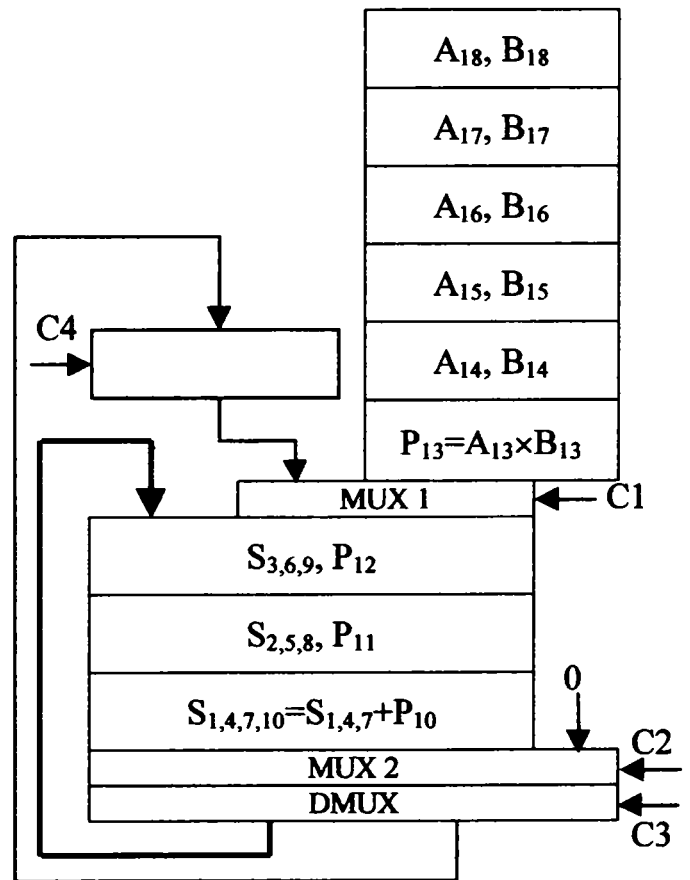
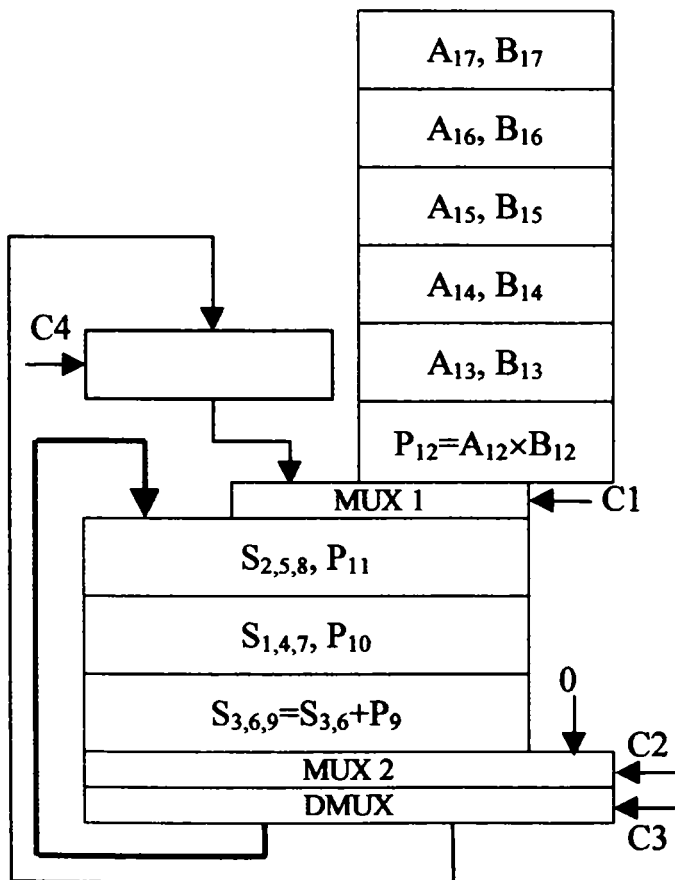
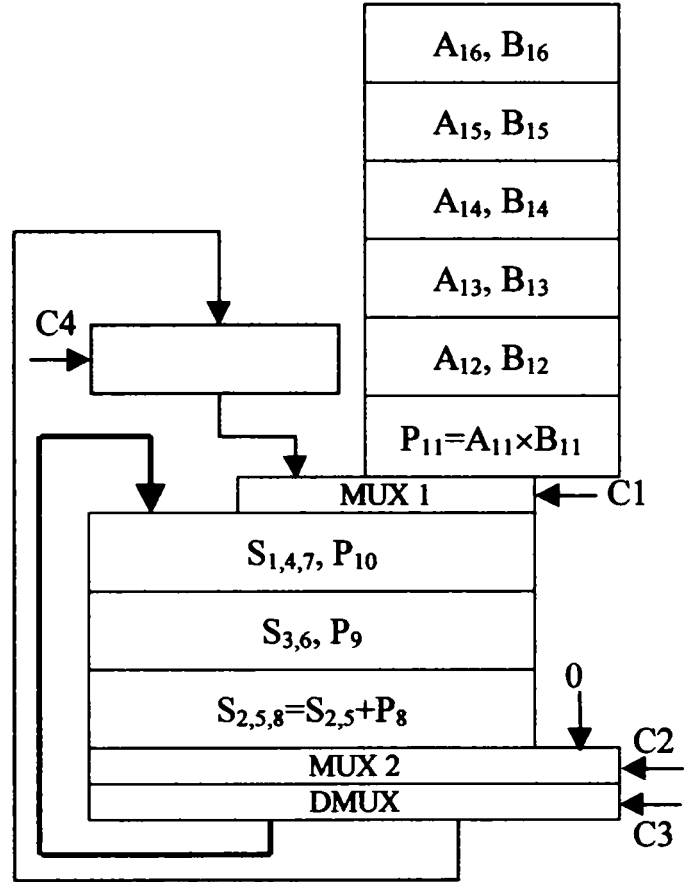
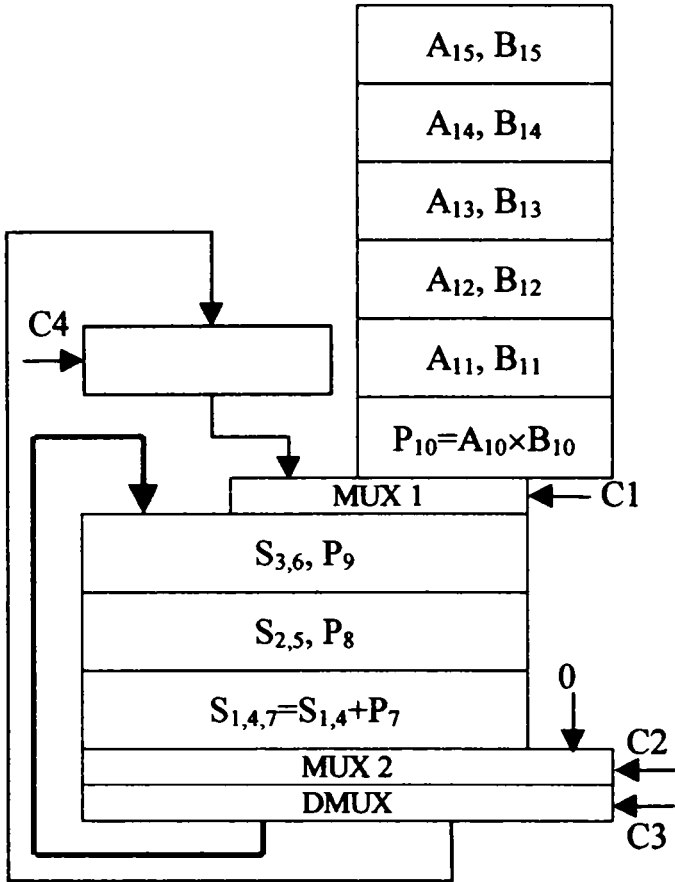
Blocurile MUX 1 și respectiv MUX 2 și DMUX se includ în nivelurile pipeline 6 din subunitatea logaritmică și respectiv 3 din subunitatea în virgulă flotantă întrucât există rezerve de 0,4 ns în fiecare din ele. În acest exemplu procesorul aritmetic calculează următoarea expresie:

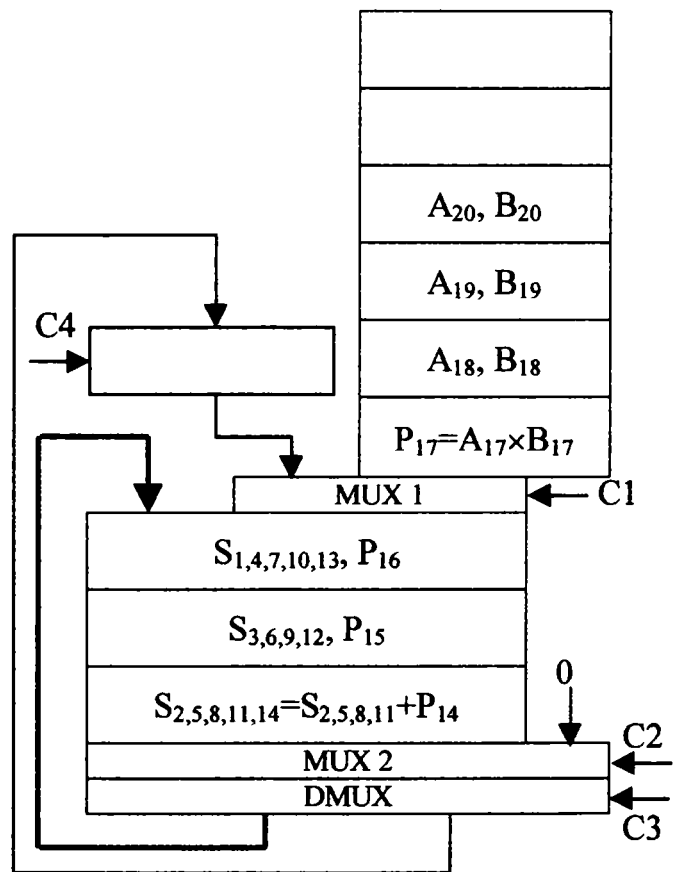
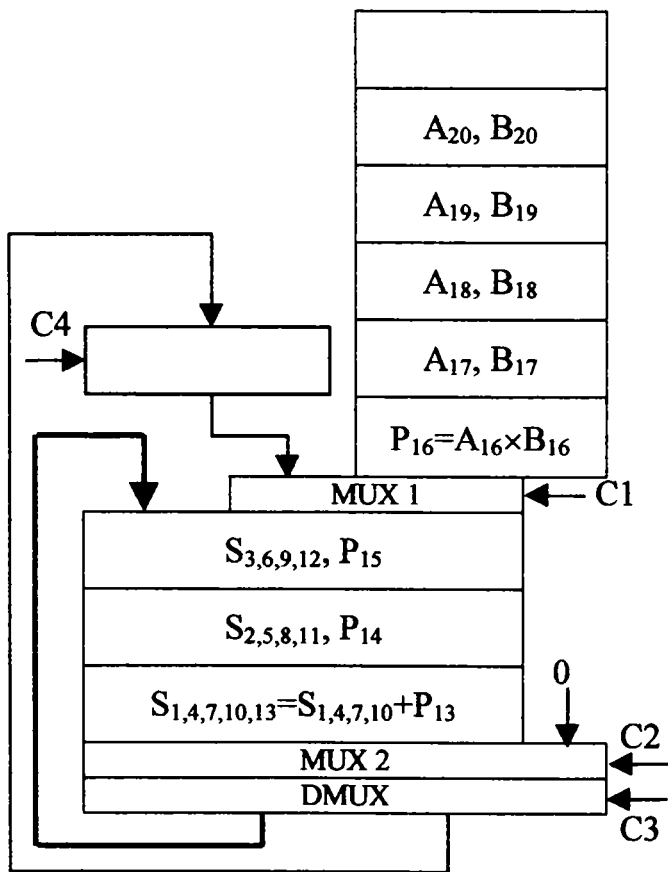
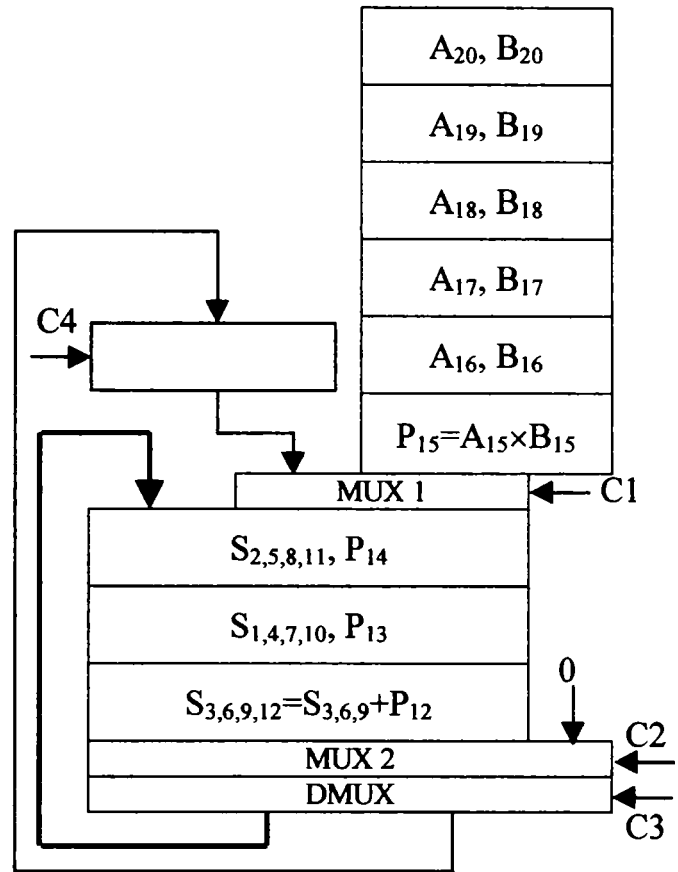
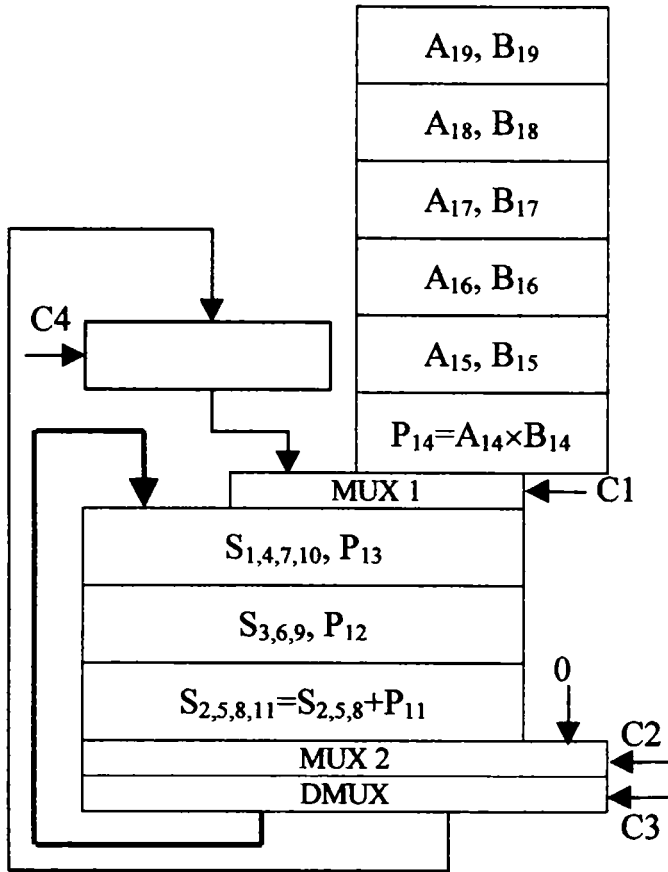
$$E = \sum_{i=1}^{20} A_i \times B_i = \sum_{i=1}^{20} P_i$$

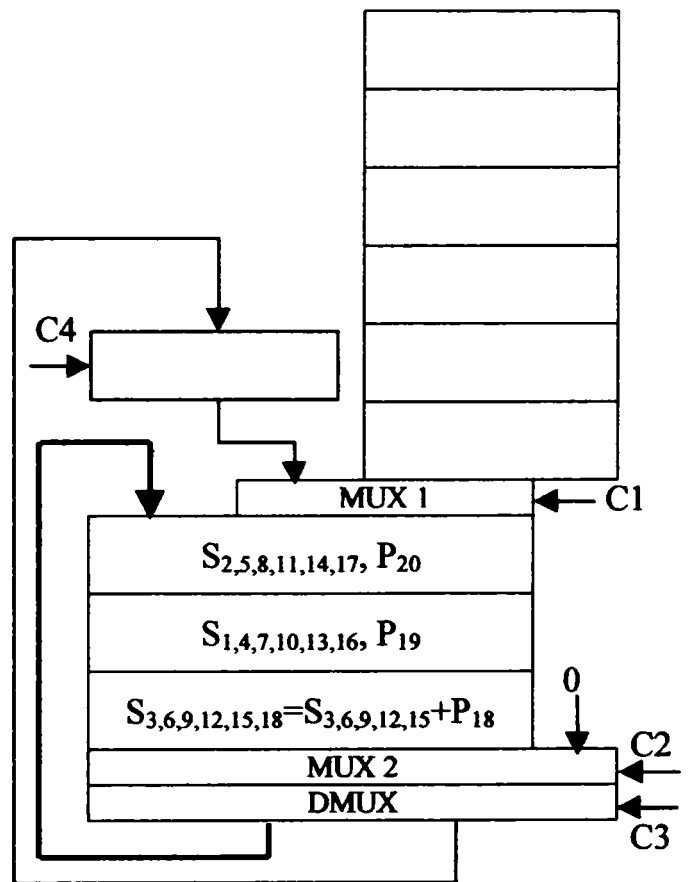
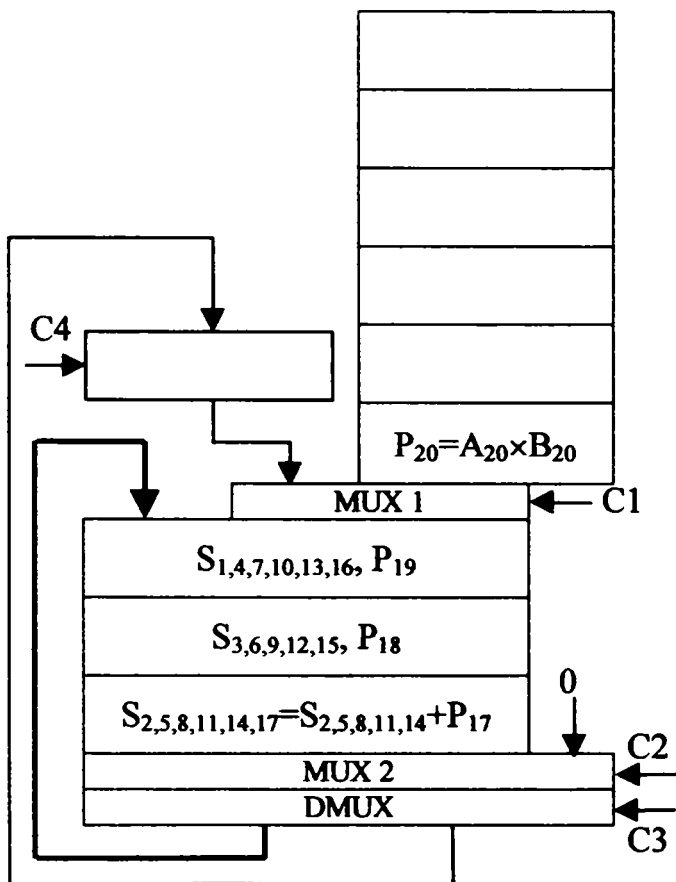
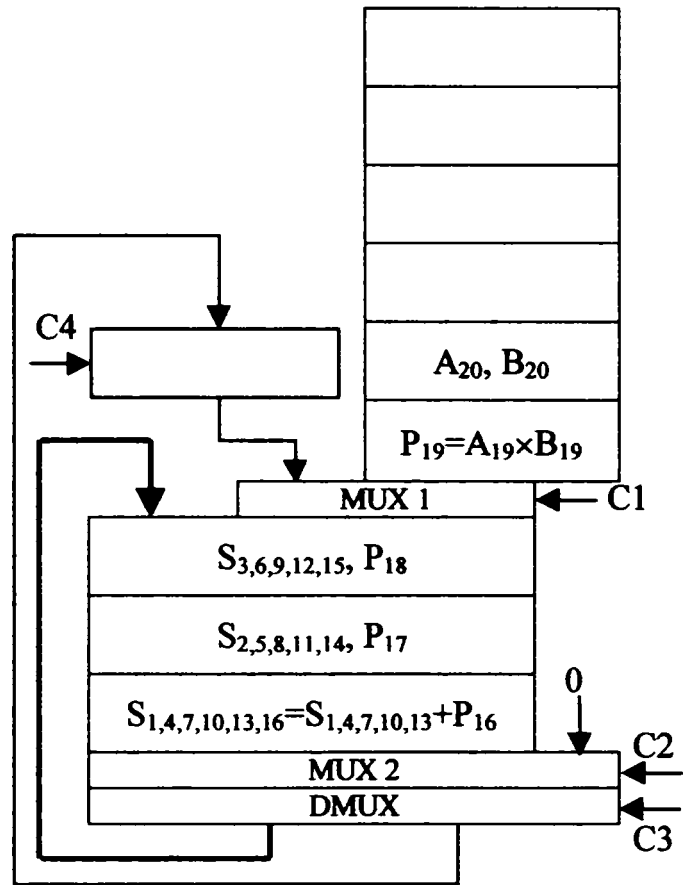
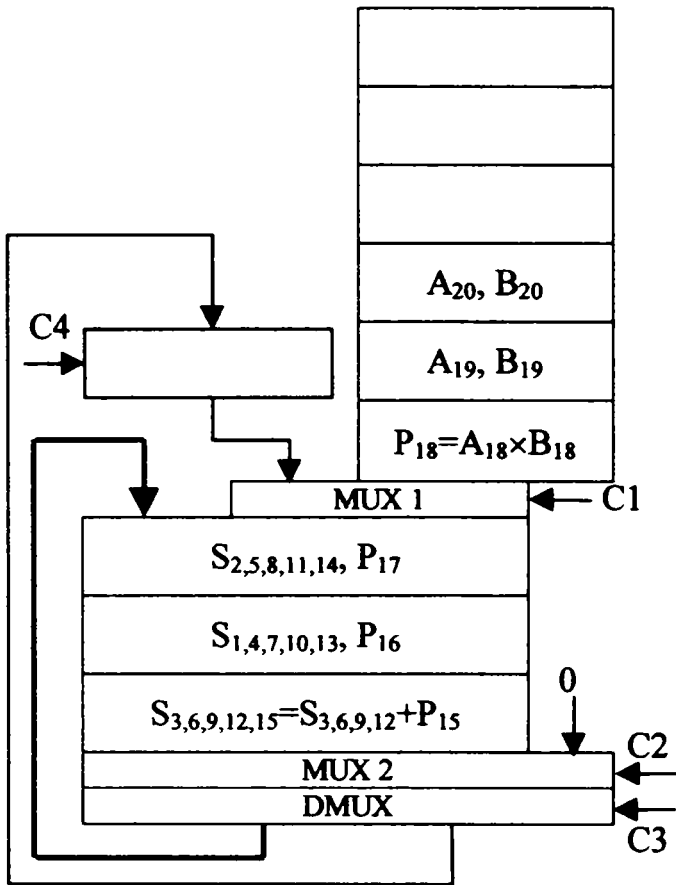


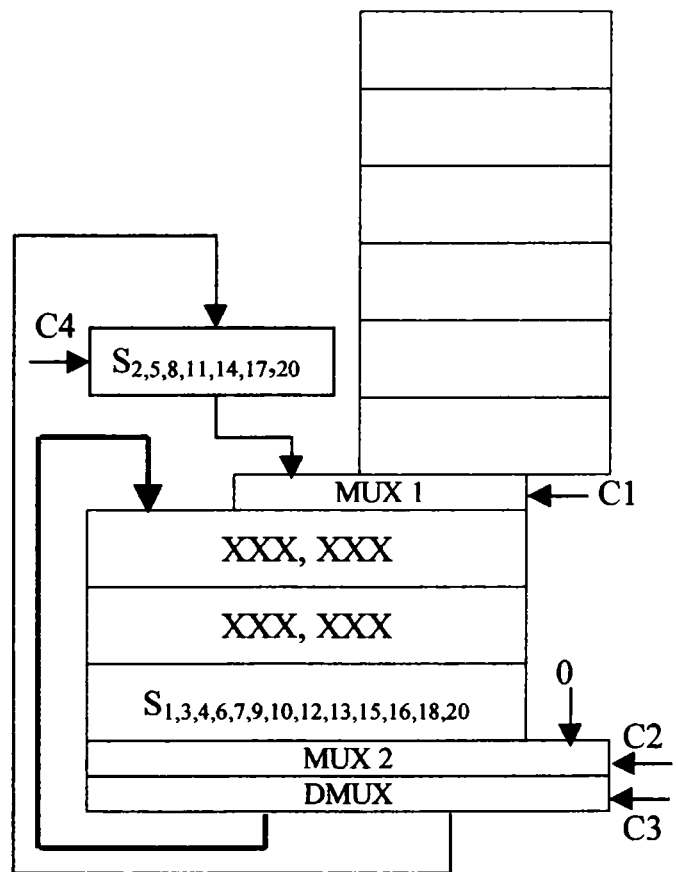
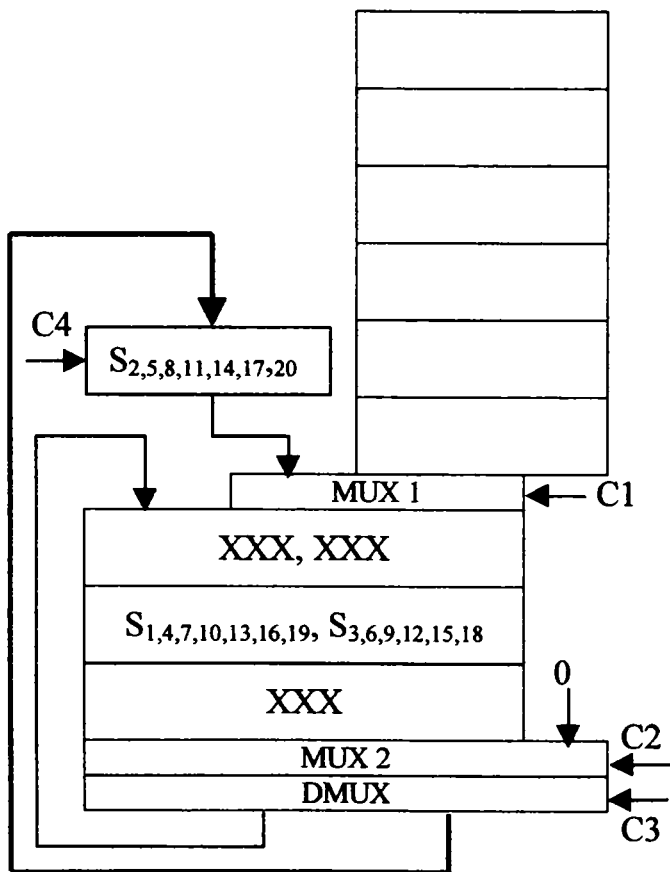
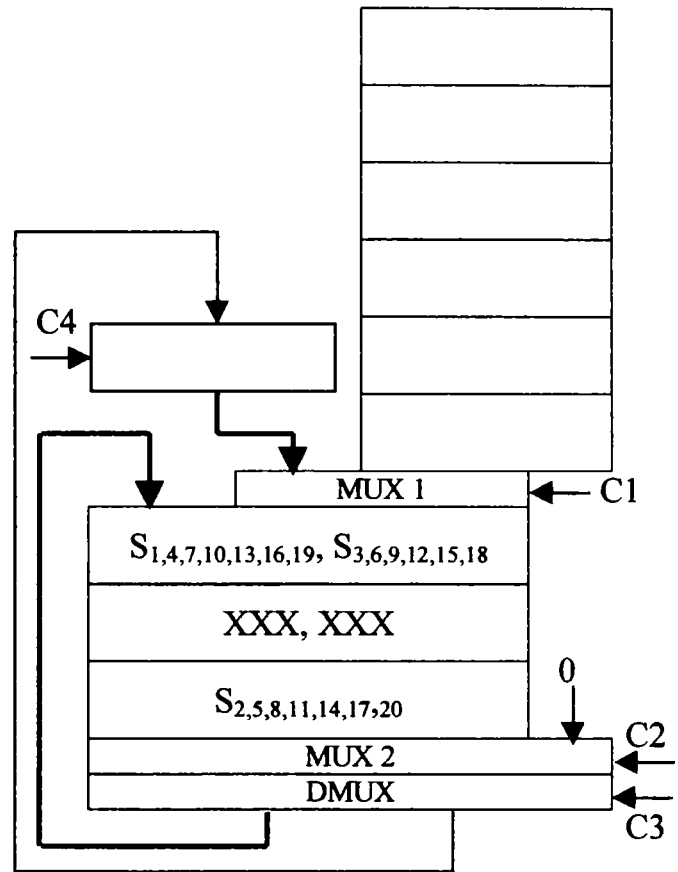
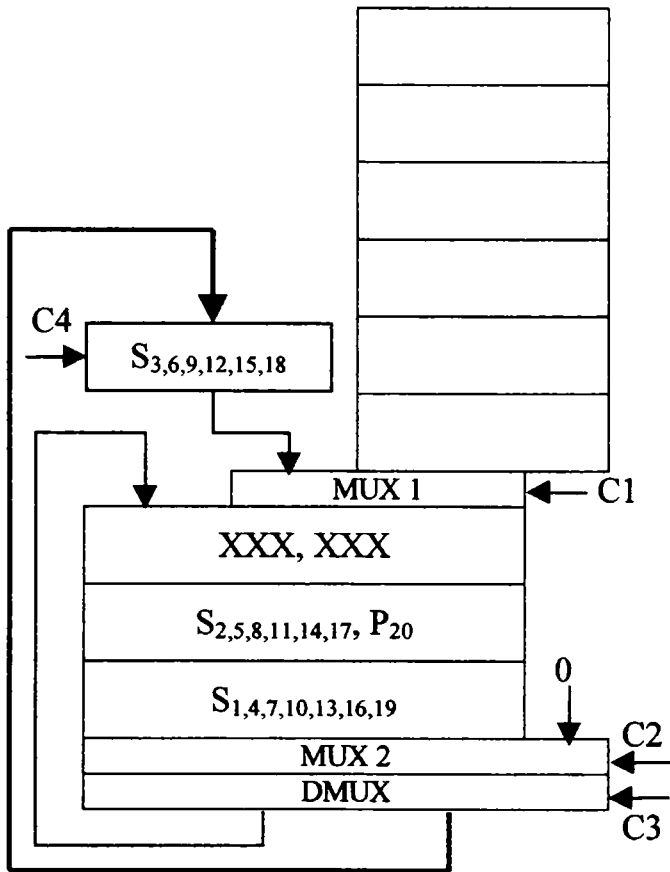
A8.1

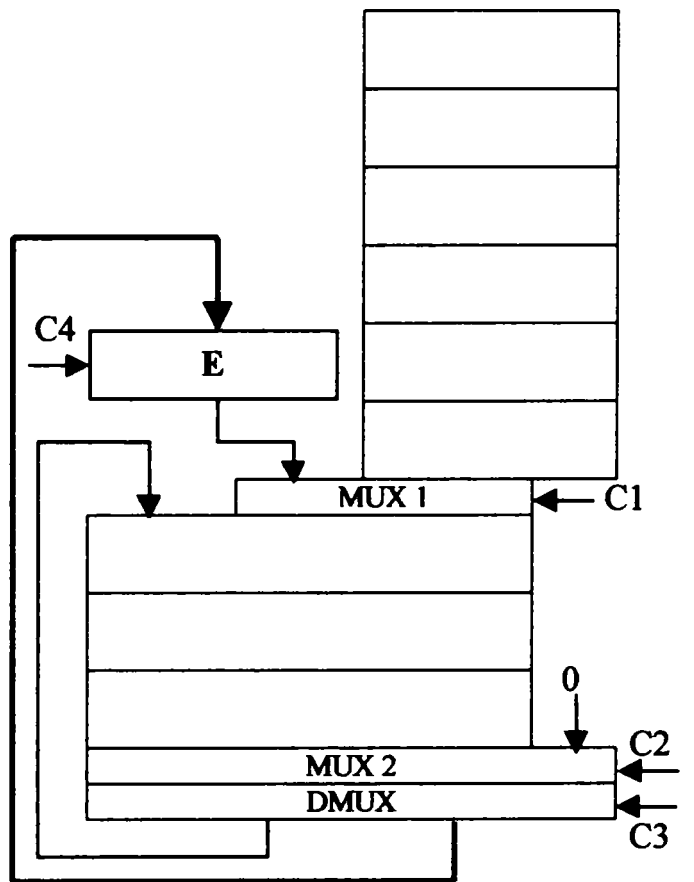
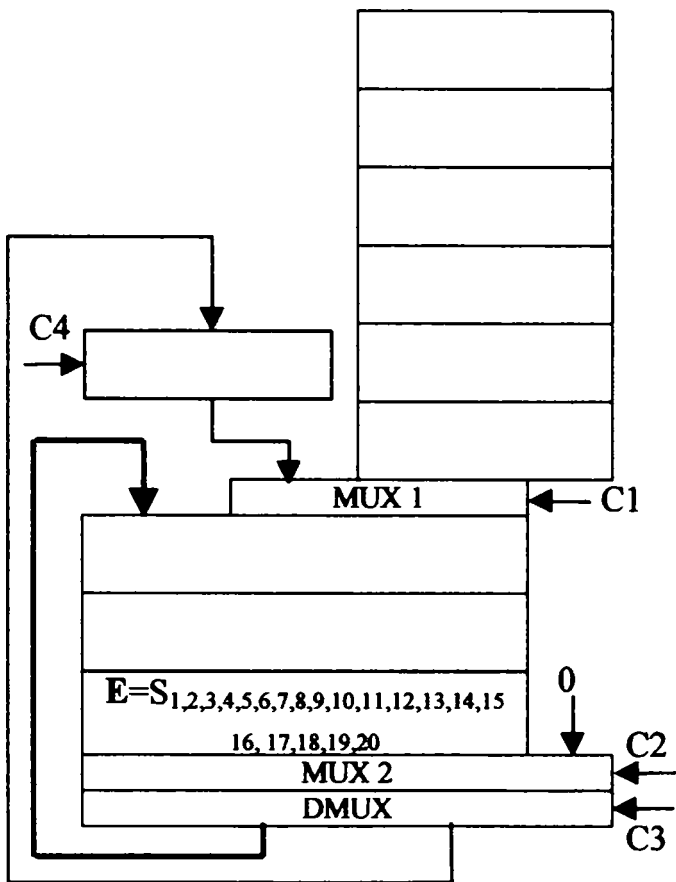
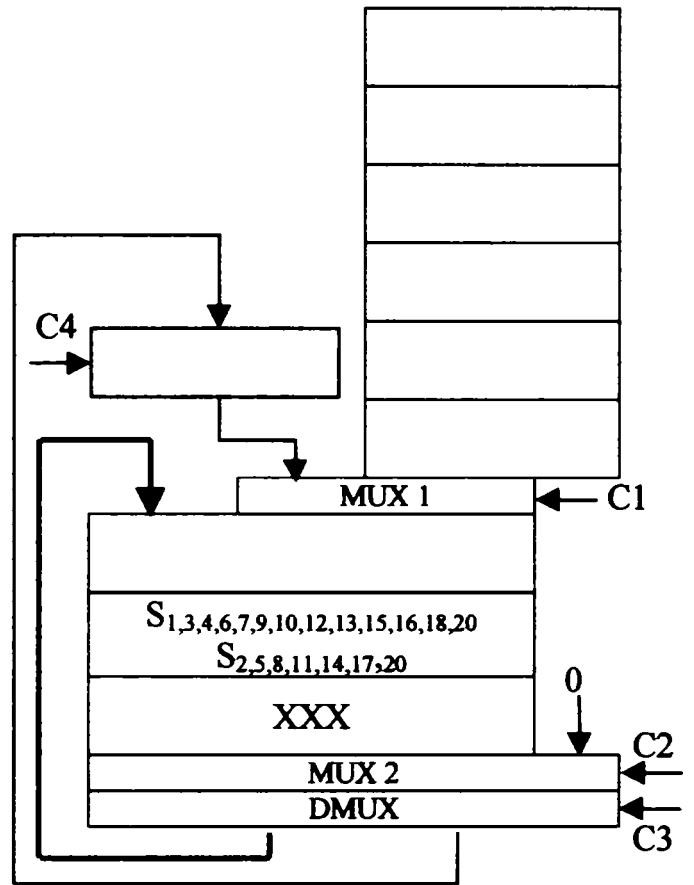
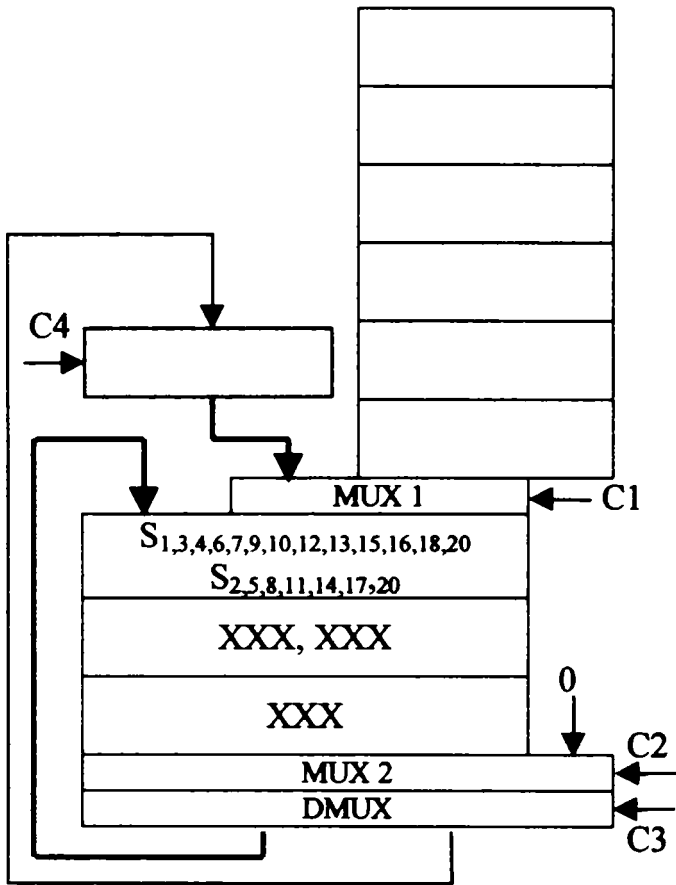












Index alfanumeric pentru circuitele prezentate în anexe

2:2 A3.9
3:2 A3.12
4:2 A3.11
A' A1.6
A' Co=0 A1.6
A' Co=1 A1.8
A1 A5.8
A2 A5.8
AA A3.5
ALU A3.17
ARB A3.9
B' A1.6
B' Co=0 A1.6
B' Co=1 A1.8
B'' A1.6
B'' Co=0 A1.6
B'' Co=1 A1.8
BB A3.4
Bloc Cprs A3.21
Bloc Logic A6.9
Bloc Ramura A-Arbore Wallace A3.1
BS A5.24
BS - Circuit "barrel shifter" pt. normalizare A7.6
BS - Circuit "barrel shifter" pt. normalizare - detaliu A7.7
BS A (pt. varianta 2) A5.15
BS B (pt. varianta 2) A5.15
BS_A A5.6
BS_A(B) detaliu A5.7
BS_B A5.6
BS_Hi A5.17
BS_Hi detaliu A5.18
Bus/6 A7.2
CC A3.6
CC1 A3.6
CC2 A3.6
CC3 A3.7
CC4 A3.7
CC5 A3.7
CC6 A3.8
CC7 A3.8
CC8 A3.8

CIRC. DEPL. MANT.-varianta 1	A5.1
CIRC. DEPL. MANT.-varianta 2	A5.11
CIRC. DEPL. MANT.-varianta 3	A5.21
CIRCUIT DE NORMALIZARE	A7.1
ExpA-ExpB	A5.3
ExpA-ExpB (pt. variantele 2 și 3)	A5.12
ExpB-ExpA	A5.4
ExpB-ExpA (pt. variantele 2 și 3)	A5.12
GPP	A3.3
Inv.B1	A3.19
Inv.B2	A3.20
Inv/NeinvA	A6.2
Inv/NeinvB	A6.2
Inv/NeinvS1	A6.7
Inversor1	A1.2
Inversor2	A1.3
InversorA	A5.2
InversorB	A5.2
LG/CPRS1	A3.11
LG/CPRS2	A3.12
LG/CPRS3	A3.13
LG/CPRS4	A3.14
LG/CPRS5	A3.15
LG/CPRS6	A3.16
LG/CSA	A3.10
LG/INV	A3.2
LG/MUX1	A3.2
LG/MUX2	A3.4
M	A5.8
M_X	A7.3
M \bar{B}	A5.2
MBA	A5.14
MI	A7.4
Mlt	A5.2
MUX	A1.2
Mux2-1	A1.8
MuxB	A1.8
MuxB	A5.8
MX	A3.9
MXI	A7.2
NZ_16	A7.4
NZ_24	A7.5
NZ_4	A7.2

NZ_8 A7.3
 NZ_8' A7.3
 RM A3.9
 RMF A3.9
 Secțiunea simulată din ALU A3.18
 Sel_A A5.9
 Sel_B A5.10
 SEL_BS A5.16
 SEL_C-da A5.14
 SEL_C-da A5.22
 SEL_MAN A5.23
 SEL_ManA A5.19
 SEL_ManB A5.20
 SEL4 A5.5
 SEL5 A5.5
 Selector A1.3
 Selector (sumator/scăzător-nivel pipeline 2) A6.8
 Selector1 A1.7
 Selector2 A1.7
 SELECTOR2-Sum1(2)_25b A6.6
 Selector3 A1.7
 SELECTOR4 A3.25
 SelExp A5.2
 SelS A5.14
 SelT A5.14
 Sum1_36b A3.22
 Sum2_36b A3.23
 Sum4_Cin=0 A3.24
 Sum4_Cin=1 A3.24
 SUM8_Cin=0 A1.5
 SUM8_Cin=1 A1.5
 SUM9_Cin=0 A6.5
 SUM9_Cin=1 A6.5
 Sumator/scăzător 32b A1.1
 SUMATOR/SCĂZĂTOR-NIVEL PIPELINE A6.1
 Sumator1 A1.4
 Sumator1_25b (varianta 1) A6.3
 Sumator1_25b (varianta 2) A6.10
 Sumator1_25b (varianta 2) detaliu A6.12
 Sumator2 A1.4
 Sumator2_25b (varianta 1) A6.4
 Sumator2_25b (varianta 2) A6.11
 Sumator2_25b (varianta 2) detaliu A6.12