

REVISTA DE ANALIZĂ NUMERICĂ ȘI TEORIA APROXIMATIEI
Volumul 3, Fascicola 1, 1974, pp. 35–51

**UN ALGORITM PENTRU GENERAREA
CUVINTELOR LIMBII ROMÂNE**

de

MINERVA BOCSA

(Timișoara)

§ 1.

Se știe că în numeroase limbi naturale cuvintele pot fi clasificate în flexibile și neflexibile. *Cuvîntul neflexibil* are o singură realizare pe planul expresiei. *Cuvîntul flexibil* poate fi conceput ca o mulțime de elemente — numite forme flexionare ale cuvîntului unitate lexicală — între care există raportul pe plan morfematic și semantic. În numeroase împrejurări cînd se enumeră cuvintele unei limbi sau porțiuni de limbă, în dicționare, liste de cuvinte, analize de texte, etc., se au în vedere doar *cuvintele unități lexicale*. Se presupune că cel care le utilizează poate reconstitui întreaga mulțime de forme flexionare, adică întreaga *paradigmă* a fiecărui cuvînt.

În prelucrarea informației lingvistice cu ajutorul calculatorului, traducere automată, rezumare automată, stabilirea caracteristicilor canticative ale limbii s.a. e nevoie să se realizeze *generarea cuvintelor și analiza acestora* prin program. *Generarea automată a cuvintelor* (sinteza) înseamnă producerea automată a întregii paradigmă, pornind de la cuvîntul unitate lexicală, ca dată. *Analiza automată a cuvintelor* înseamnă stabilirea cuvîntului unitate lexicală, adică apartenența la o anumită paradigmă, pornind de la un cuvînt oarecare al limbii, ca dată.

Ne-am propus să rezolvăm generarea cuvintelor limbii române cu ajutoul calculatorului IRIS 50. Ne-am mărginit la un număr de cca. 2000 cuvinte regulate, substantive, adjective și verbe, majoritatea cuprinse în liste de cuvinte din „*Morfologia structurală a limbii române*” de Valeria Guju Romalo, [6], avînd în vedere doar *formele sintetice* (compuse dintr-un singur cuvînt). Am ales limbajul ASSIRIS deoarece dispune de o listă largă de instrucții adecvate problemei. Prezentăm în cele ce urmează algoritmul pentru generarea cuvintelor *flexiunii nominale și verbale* din limba română.

§ 2.

Cuvîntul flexibil se obține prin juxtapunerea a două elemente: o parte constantă, numită *radical* (rădăcină, temă) și o parte variabilă, numită *flectiv* (terminație, modifier). Flectivele posibile se grupează în clase,

dind naștere la *clasele flexionare* ale limbii. Avem substantive de declinarea I-a, a II-a, ... verbe de conjugarea I-a, a II-a ... Astfel substantivul MAMĂ din declinarea I-a are ca radical MAM- iar seria de flective este -Ă, -E, -A, -EI, -ELE, -ELOR. Paradigma substantivului MAMĂ este formată din cuvintele distințe MAM-Ă, MAM-E, MAM-A, MAM-EI, MAM-ELE, MAM-ELOR.

Pentru realizarea cuvintelor ar fi suficient să separăm radicalul și să atașăm succesiv seria de flective. Un asemenea algoritm a fost utilizat de ERICA NISTOR DOMOKOS [11], în experimentele sale de traducere automată din limba engleză în limba română pe calculatorul MECIPT-1, în 1962 și anii următori. O variantă a sa mai amplificată a fost la baza sintezei substantivelor rusești cu calculatorul DACICC-1, imaginată de Paul Schveiger și LIVIU NEGRESCU, 1966, [12].

Limba română însă are multe neregularități care nu permit generalizarea acestui algoritm. Există cuvinte care nu au un radical constant. Cuvântul MASĂ are radicalele MAS-(Ă) și MES-(E), cuvântul MĂR admite radicalele MĂR-(Ø)* și MER-(E), cuvântul FLOARE are radicalele FLOAR-(E) și FLOR-(I), cuvântul CAL are radicalele CAL-(Ø, -UL, -ULUI) și CA-(I, -II, -ILOR). Aceste forme diferite ale radicalului se numesc *alomorfe* iar modificările de tip A-E, Ă-E, A-Ø, L-Ø, etc. se numesc *alternanțe fonetice*.

S-a stabilit că în limba română radicalele substantivale au maximum două alomorfe, cele adjecтивale au trei, FRUMOS-Ø, FRUMOȘ-I, FRUMOAS-E) iar verbele ajung la cinci și șase alomorfe, fără totuși să prezinte ceea ce se numește o variație aberantă. (VED-EA, VĂD-Ø, VEZ-I, VAD-Ă, VĂZ-IND; RAD-E, RAZ-I, RĂD-EAM, RĂ-SEI, RA-SE, RĂZ-IND).

Alternanțele fonetice (vocalice și consonantice) se produc foarte frecvent și nu sunt circumscrise decât aproximativ de condiționări fonetice. (CAS-Ă, apropiat de MAS-Ă și din aceeași declinare are totuși un singur radical, SACAL- din aceeași declinare cu CAL-Ø nu-l pierde pe L, în cursul paradigmăi s.a.m.d.). Nici segmentarea cuvântului în alt mod — problemă controversată — nu detașează net partea fixă din cuvânt de partea variabilă. Singurul mod de descriere utilizabil în problema noastră este *stabilirea unui inventar complet* de astfel de alternanțe și stabilirea repartiției valorilor în paradigmă.

Considerăm că o asemenea descriere se face în cel mai economic mod utilizând litere variabile, noțiune concepută de GR. C. MOISIL în cercetările sale legate de problemele traducerii automate [9], [10]. Literă variabilă T/T din BĂRBAT-Ø are valoarea (*realizarea*) T în (a) 1-Nom., sing., neart., 2-Gen., sing., neart., 3-Dat., sing., neart., 4-Ac., sing., neart., etc. și valoarea T în (b) 1-Nom. pl., neart., 2-Gen., pl., art., etc. Astfel întîlnim literă variabilă A/A din PACE (PA/ĀCE), care dă PACE, PĀCI; literele variabile Z/J și A/Ø din VITEAZ (VITEA/ØZ/J) care dă VITEAZ, VITEJI; literele variabile ř/S, T/C, Ø/A din CREȘTE (CREØ/AS/ST/CE) care dă CRESC, CREŠTI, CREASCĂ s.a.m.d.

* Simbolul vid, lanțul vid, numit și zero.

Regăsim literele variabile în [6] unde au servit la descrierea structurală a substantivului, adjecтивului și verbului.

Sinteza automată a cuvintelor trebuie să prevadă două aspecte:

- I. o codificare a întregii informații necesare flexionării;
- II. un program de decodificare și de construire a paradigmăi.

§ 3.

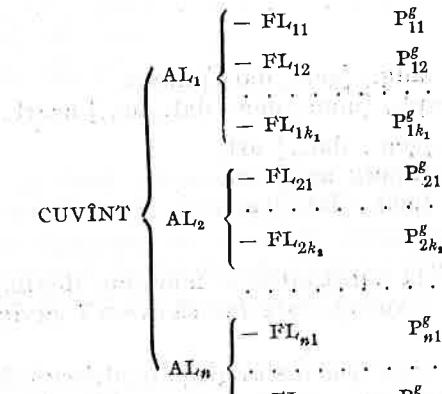
Codificarea informației se face ținând cont de fazele operației:

- A) cuvântul de pe cartela perforată se introduce în memoria centrală.
- B) programul prelucrează cuvântul dat,
- C) se tipărește la imprimantă paradaima obținută.

Se știe că la calculatorul IRIS 50 cititorul de cartele preia informația de intrare de pe cartelele perforate într-un anumit cod (Hollerith) și o introduce în memoria centrală unde se codifică în EBCDIC.

Codificarea literelor variabile trebuie să conțină următoarele elemente:

- a) *realizările* (valorile efective) ale literei variabile în toate alomorfele radicalului;
 - b) repartizările radicalului în paradigmă. Astfel codificarea trebuie să permită construirea paradigmăi după schema (*).
- Aici \overline{AL}_i cu $i = 1, n$ sunt cele n alomorfe ale radicalului, FL_{ij} cu $j = 1, k_i$ sunt cele k_i flective ce se atachează radicalului AL_i , iar P_{ij}^g sunt pozițiile



Schema de formare a paradigmăi. (*)

din paradigmă ce se realizează prin lanțul AL_i, FL_{ij} . Numărul g arată gruparea pozițiilor P_{ij}^g în cazul concret al fiecărui cuvânt în parte.

La *substantive* gruparea elementelor P_{ij}^g e simplă, reducindu-se la două (trei) situații standard, substantive masculine și feminine, excepțiile fiind rare. Numărul total de forme flexionare de la substantive este la noi 16:

număr — sing. sau plur. —, *caz* — nom. gen., dat., ac., — și *determinarea* — neart., art. (nu am considerat vocativul). Dar numeroasele omonimii reduc paradigmă la 6 (7) cuvinte diferite. Cele două alomorfe posibile se repartizează în grupările (1) și (2) după cum urmează :

(1) Substantive masculine si neutre

AL_1 { CM1: sing., [nom., gen., dat., ac.,] neart. (BÄIAT-Ø)
 CM2: sing., [nom., ac.,] artic. (BÄIAT-UL)
 CM3: sing., [gen., dat.,] artic. (BÄIAT-ULL)

| | | | | | | | | | | |
|------------------------------------|--|------------------------------------|--------|-----------|------------------------|------|------------|-------------------------|------|---------------|
| AL_2 | <table border="0"> <tr> <td>CM4: pl., [nom., gen., dat., ac.,]</td><td>neart.</td><td>(BÄIET-I)</td></tr> <tr> <td>CM5: pl., [nom., ac.,]</td><td>art.</td><td>(BÄIET-II)</td></tr> <tr> <td>CM6: pl., [gen., dat.,]</td><td>art.</td><td>(BÄIET-II,OR)</td></tr> </table> | CM4: pl., [nom., gen., dat., ac.,] | neart. | (BÄIET-I) | CM5: pl., [nom., ac.,] | art. | (BÄIET-II) | CM6: pl., [gen., dat.,] | art. | (BÄIET-II,OR) |
| CM4: pl., [nom., gen., dat., ac.,] | neart. | (BÄIET-I) | | | | | | | | |
| CM5: pl., [nom., ac.,] | art. | (BÄIET-II) | | | | | | | | |
| CM6: pl., [gen., dat.,] | art. | (BÄIET-II,OR) | | | | | | | | |

(2) Substantive feminine

$\text{AL}_1 \left\{ \begin{array}{l} \text{CF1: sing., [nom., ac.,] neart.} \\ \text{CF2: sing., [nom., ac.,] art.} \end{array} \right. \quad (\text{FAT-A})$

| | | |
|--------|---|--|
| AL_2 | $\left\{ \begin{array}{l} CF_3: \left\{ \begin{array}{l} CF_3': \text{sing., [gen., dat.,]} \text{ neart.} \\ CF_3'': \text{pl., [nom., gen., dat., ac.,]} \text{ neart.} \end{array} \right. \\ CF_4: \text{sing., [gen., dat.,]} \text{ art.} \\ CF_5: \text{pl., [nom., ac.,]} \text{ art.} \\ CF_6: \text{pl., [gen., dat.,]} \text{ art.} \end{array} \right.$ | (FET-E) (FET-EL) (FET-ELB) (FET-ELOR) |
|--------|---|--|

În cazuri izolate la substantivele feminine distingem cr³ diferit de CF^{3''} (cuvintele SARE, CARNE), care fac să avem 7 cuvinte diferite în paradigmă.

Cuvîntul VRABIE), cere însă o altă grupare, diferită de (1) și (2), deoarece CF4 are același alomorf cu CF1 (VRABIEI). (Circulă paralel și VRĂBIEI, întocmai ca și ARAMEI și ARĂMII)

Din analiza celor circa 1.300 *substantive* am desprins în total două grupări obișnuite, (1) și (2) și încă două grupări *aberante* (3) și (4). Aceasta ne-a condus la codificarea și programul expus în [1].

Adjectivele studiate (circa 250) au la noi 32 de forme flexionare, masculine și feminine, respectând în principiu aceleași omonimii ca [substantivele și cu un număr mic de grupări ale elementelor paradigmăi, care ne-a condus la codificarea și programul expus în [2].

Verbele studiate (circa 450) prezintă însă o bogată variație a acestor grupări. [6]. De exemplu la verbe cu *două alomorfe* ale radicalului avem grupările:

| | | | |
|-----|--------|--|--|
| | AL_1 | { infinitiv ind. prez. I. sing. ind. prez. III. sing. (toate celelalte) | (LUPT-A) (LUPT-Ø) (LUPT-Ä) (LUPT-...) |
| (1) | AL_2 | { ind. prez. II. sing. | (LUPT-I) |
| | | - - - - - | - - - - - |
| | AL_1 | { infinitiv ind. prez. I. sing. ind. prez. III. sing. (toate celelalte) | (APÄR-A) (APÄR-Ø) (APÄR-Ä) (APÄR-...) |
| (2) | AL_2 | { ind. prez. II. sing. conj. prez. III. sing. conj. prez. III. pl. | (APER-I) (APER-E) (APER-E) |
| | | etc. | |

Pentru verbele cu două alomorfe avem grupările numerotate (1), (2), (11), la cele cu trei alomorfe avem grupările (1), ... (22) și a.m.d.

Literele variabile figurează în alomorfe conform unor configurații, care de asemenea diferă mult. Astfel literele E/I și T/Ț din verbul A PREZENT-A și literele A/Ӑ și C/Ҫ din verbul A FAC-E (ambele cu trei alomorfe ale radicalului) dau realizările:

| | | |
|----------------|----------------|----------------|
| AL_{t_1} (a) | AL_{t_2} (b) | AL_{t_3} (c) |
| PREZENT-A | PREZINT-Ø | PREZINT- |
| FAC-E | FÄC-EAM | FÄØ-Ø |

Ayem — notînd alomorfele cu (a), (b), (c) — literele

E/I în configurația (a)/(b)(c)
 T/T în configurația (a)(b)/(c)
 A/Ā în configurația (a)/(b)(c)
 C/Q în configurația (a)(b)/(c).

Literele variabile de la trei alomorfe prezintă și alte configurații (a)(c)/(b), (a)/(b)/(c), la patru alomorfe avem configurațiile (a)(b)/(c)(d), etc. Numerotarea alomorfelor $AL_1, AL_2 \dots AL_n$, respectiv (a), (b), (c), ..., se face după *temeiuri lingvistice* să incit nu putem — renumerotindu-le — să considerăm ca egale configurațiile (a)(b)/(c) și (a)(c)/(b).

În vederea codificării tuturor informațiilor pentru substantive, adjective și verbe, am exploatat posibilitățile calculatorului IRIS 50 care admite o unitate minimă adresabilă de informație *octetul* (opt cifre binare) format din cvarțetul stâng (ponderea tare, poids fort) și cvarțetul drept (ponderea slabă, poids faible). Am alcătuit un inventar al *literelor variabile generalizat al limbii române*. Acest alfabet conține literele:

A, Ă, B, C, D, E, F, G, H, I, Ī, J, K, L, M, N, O, P, Q, R, S, ř, T, ū, V, W, X, Y, Z, blanc*) precum și cele 82 litere variabile din care dăm mai jos cîteva:

| Nr. | Lit. | Nr. atom. | Config. | Cuvînt |
|-----|-------|-----------|---------------------|------------------------------------|
| 1 | A/Ā | 2 | (a)/(b) | CART-E/CĂRT-I |
| 9 | Ø/A | 2 | (a)/(b) | MOTOR-Ø/MOTOAR-E |
| 23 | N/Ø | 2 | (a)/(b) | SPUN-E/SPU-I |
| 25 | A/A | 3 | (a)/(b)/(c) | BAT-E, BAT-I/BĂT-EAM |
| 49 | U/O | 3 | (a)/(b)/(c) | RUG-A/ROG-Ø, ROAG-Ā |
| 54 | Ā/E/A | 3 | (a)/(b)/(c) | SPĂL-A/SPEL-I/SPAL-Ā |
| 58 | C/P | 4 | (a)/(b)/(c)/(d) | COC-Ø, COAC-E/COP-T, COAP-SERĂ |
| 79 | D/Z/Ø | 5 | (a)/(b)/(c)/(d)/(e) | ROAD-E, ROD-Ø/ROZ-I/RO-SEI, ROA-SE |

Alfabetul obținut, inclusiv numărul alomorfelor radicalului și configurațiile, s-a codificat cu ajutorul codului CGCLR (Codul general al cuvîntelor limbii române), ce figurează în tabelul 1.

Se observă că literele fixe au ponderea tare E și F, literele cu două alomorfe ale radicalului au ponderea tare C și D literele ce comută cu zero au ponderile tari 0, 1, 2, 3, 4, s.a.m.d. În general tabelul nr. 1 conține o mare regularitate privind decodificarea, deoarece ponderea tare (coloana) arată numărul alomorfelor și în multe cazuri configurația, iar ponderea slabă (linia) arată valorile reale ale literelor variabile. Au existat însă litere a căror decodificare nu e aşa regulată. Ele ocupă zona dublu încadrată de la coloanele 2, 3, 4, 5, 6, 8, C.

Alfabetul nostru poate fi lărgit dacă se vor analiza în acest mod toate cuvîntele limbii române, literele variabile noi fiind plasate în spațiile libere din tabelul 1.

În tabelul 2 dăm elementele principale ale codului EBCDIC.

*) Literele Ă, Ī, ř, ū, T, care nu figurează la calculatorul IRIS-50 sunt codificate cu semnele %,], \$, &.

§ 4.

In memoria internă a calculatorului cuvîntul va fi reprezentat în CGCLR. Pe cartelă îl vom codifica astfel:

a) Literele fixe le vom reprezenta printr-un caracter corespunzător literei din alfabetul obișnuit;

b) Literele cu variație regulată vor fi codificate prin trei caractere: un *prefix* format de două cifre zecimale, indicînd coloana din tabelul 1, urmat de una din valorile *realizărilor*. De ex. litera nr. 1 (A/Ā) este 13A, urmat de una din valorile *realizărilor*. De ex. litera nr. 1 (A/Ā) este 13A, litera nr. 9 (Ø/A) este 00A, litera nr. 25 (A/A) este 11A, litera 49 (U/O) este 08U s.a.m.d.

c) Literele cu variație neregulată vor fi codificate cu patru caractere: *prefixul* format de două cifre zecimale ca la b), o cifră zecimală *separator*, apoi o literă, una din valorile *realizărilor*. De ex. litera nr. 32 (Ø/A) este apoi o literă, una din valorile *realizărilor*. De ex. litera nr. 32 (Ø/A) este 020A, litera nr. 63 (A/Ø) este 030A, litera nr. 64 (A/Ø) este 031A, litera nr. 58 (C/P) este 062C s.a.m.d.

Pe fiecare cartelă figurează un singur cuvînt în felul următor:

- col. 1—25: cuvîntul cu literele codificate ca la a), b), c).
- col. 26—27: LCV, lungimea cuvîntului de pe cartelă, două cif. zec.
- col. 28—29: LRD, lungirea radicalului din forma internă, două cif. zec.
- col. 30: NAL, numărul alomorfelor radicalului, o cif. zec.
- col. 31—32: GRP, grupările alomorfelor radicalului, două cif. zec.
- col. 33: PV, partea de vorbire care poate fi numai S (subst.), A (adj.) sau V (verb), o literă.
- col. 34: GT, gen pentru substantive, blanc pentru adjective și tip pentru verbe, o literă sau blanc.
- col. 35—36: CLF, clasa flexionară, arătînd declinarea sau conjugarea, două cif. zec.
- col. 37—80: neutilizat.

De exemplu avem cuvîntele TARĂ, SPĂLA, codificate pe cartelă astfel:

| | |
|----------------------------------|--|
| col. 1, 2, 3, 4, 5, 6, 7, 8, ... | 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36 |
| T 1 3 A R Ā | 0 6 0 3 2 0 2 S F 1 3 |
| S P 0 8 4 A L A | 0 8 0 4 3 0 4 V A 0 1 |

§ 5.

Algoritm propus de noi pentru realizarea fazelor A, B, C, are etapele:

- I_A — Codifiarea informației de intrare (de pe cartelă) din codul EBCDIC în CGCLR..
- II_B — Formarea alomorfelor radicalului.
- III_B — Atașarea flectivelor
- III_C — decodificarea din CGCLR în codul EBCDIC și tipărirea informației de ieșire.

Tabloul 7

| | | Codul CGCLR | | | | | | | | | | | | | | | |
|-------|-----------------------|---------------|-----------------------|-----------------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
| | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F |
| slabă | 0000 | 0001 | 0010 | 0011 | 0100 | 0101 | 0110 | 0111 | 1000 | 1001 | 1010 | 1011 | 1100 | 1101 | 1110 | 1111 | |
| 0 | \emptyset/A | A/\emptyset | \emptyset/\emptyset | \emptyset/\emptyset | A/\emptyset | |
| 0000 | | | | | | | | | | | | | | | | | |
| 1 | \emptyset/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset |
| 0001 | | | | | | | | | | | | | | | | | |
| 2 | \emptyset/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset |
| 0010 | | | | | | | | | | | | | | | | | |
| 3 | \emptyset/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset |
| 0011 | | | | | | | | | | | | | | | | | |
| 4 | \emptyset/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset |
| 5 | \emptyset/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset |
| 0100 | | | | | | | | | | | | | | | | | |
| 6 | \emptyset/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset |
| 0110 | | | | | | | | | | | | | | | | | |
| 7 | \emptyset/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset | A/\emptyset |
| 0111 | | | | | | | | | | | | | | | | | |

8

| | | X | | | | | | | | | | | | | | | |
|-----------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|
| | | | | | | | | | | | | | | | | | |
| 8 | | | | | | | | | | | | | | | | | |
| 1000 | | | | | | | | | | | | | | | | | |
| 9 | G/\emptyset |
| 1001 | | | | | | | | | | | | | | | | | |
| A | C/\emptyset |
| 1010 | | | | | | | | | | | | | | | | | |
| B | B/\emptyset |
| 1011 | | | | | | | | | | | | | | | | | |
| C | N/\emptyset |
| 1100 | | | | | | | | | | | | | | | | | |
| D | $D/z/\emptyset$ |
| 1101 | | | | | | | | | | | | | | | | | |
| E | \emptyset/H |
| 1110 | | | | | | | | | | | | | | | | | |
| F | L/\emptyset |
| 1111 | | | | | | | | | | | | | | | | | |
| Config. | a/b |
| | Comutare | cum | zero | | | | | | | | | | | | | | |
| Nr. alom. | 2 alomorfe | 3 alom. | 4 alom. | Cazuri speciale | 4 alomorfe | 3 alomorfe | 2 alom. | | | | | | | | | | |
| Prefix | 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 | 10 | 11 | 12 | 13 | 14 | 15 | |

Tabelul 2.
Codul EBCDIC

| Pondere tare slabă | 0 0000 | 1 0001 | 2 0010 | 3 0011 | 4 0100 | 5 0101 | 6 0110 | 7 0111 | 8 1000 | 9 1001 | A 1010 | B 1011 | C 1100 | D 1101 | E 1110 | F 1111 |
|-----------------------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| | | | | | | | | | | | | | | | | |
| 0 0000 | | | | | | | | | | | | | | | | 0 |
| 1 0001 | | | | | | | | | | | | | | | | 1 |
| 2 0010 | | | | | | | | | | | | | | | | 2 |
| 3 0011 | | | | | | | | | | | | | | | | 3 |
| 4 0100 | | | | | | | | | | | | | | | | 4 |
| 5 0101 | | | | | | | | | | | | | | | | 5 |
| 6 0110 | | | | | | | | | | | | | | | | 6 |
| 7 0111 | | | | | | | | | | | | | | | | 7 |
| 8 1000 | | | | | | | | | | | | | | | | 8 |
| 9 1001 | | | | | | | | | | | | | | | | 9 |
| A 1010 | | | | | | | | | | | | | | | | |
| B 1011 | | | | | | | | | | | | | | | | |
| C 1100 | | | | | | | | | | | | | | | | |
| D 1101 | | | | | | | | | | | | | | | | |
| E 1110 | | | | | | | | | | | | | | | | |
| F 1111 | | | | | | | | | | | | | | | | |

Tabelul 2.
Codul EBCDIC

Nucleul întregului program îl formează instrucția de traducere și testare pe care o utilizăm de la stînga spre dreapta (TRTR, Translate and test). Această instrucție prelucrează un lanț de caractere de lungime dată, executând asupra lui ambele sau una sau cealaltă din funcțiuni: traduce, adică înlocuiește fiecare caracter (octet) cu un alt caracter, conform unei tabele de corespondențe; testează, adică verifică dacă octetul are sau nu o proprietate.

Proprietatea poate fi foarte variată: octetul dat, întreg sau orice grupare de biți ai săi, selecționată prin cifrele 1 ale unui alt octet, denumit *masca de comparare*, se compară cu un alt treilea octet, numit *criteriu de comparare*. De exemplu cifrele din ponderea tare (biții 0, 1, 2, 3) se selecționează prin masca 1111 0000, iar cifrele din ponderea slabă (biții 4, 5, 6, 7) se selecționează prin masca 0000 1111. Compararea se face cu biții corespunzători din octetul „criteriu de comparare” și este de forma: caracterul dat este <, <=, =, >, >=, ≠ față de criteriul de comparare. Execuția instrucției se oprește în interiorul lanțului, cînd s-a aflat proprietatea testată, sau dacă s-a epuizat tot lanțul fără găsirea caracterului căutat.

Etapa IA a algoritmului are schema logică din figura 1.

Instrucția TRTR are aici ambele funcții. Traduce literele cuvîntului de lungime LCV (șirul de lungime L) după o tabelă de decodificare (tabelul 3). În același timp testează prezența cifrei în lanț (adică a literei variabile), verificînd dacă există un caracter cu ponderea tare egală cu F. (vezi tabelul 2). Apoi se înlocuiesc fie următoarele trei, fie patru caractere cu un singur octet, conform cu tabelul 1.

Tabelul 3

| Simbol literă | Cod EBCDIC | Cod CGCLR |
|---------------|------------|-----------|
| A | C1 | E0 |
| B | C2 | FC |
| C | C3 | EA |
| ... | ... | ... |
| M | D4 | F3 |
| N | D5 | FD |
| P | D7 | FA |
| ... | ... | ... |

După această etapă cuvintele TARĂ, FRUMOS, SPĂLA, vor fi reprezentate în calculator (cod hexa).

TA/Ā R Ā;
EC D0 F5 E1;
F R U M O Ø/A S/\$;

F2 F5 E5 F3 E6 21 A7;

E7 FA 8F FF E0

Etapa II_B are schema logică de ansamblu în figura 2, iar rutina de analiză și depunerea literelor variabile în alomorfele radicalului în figura 3.

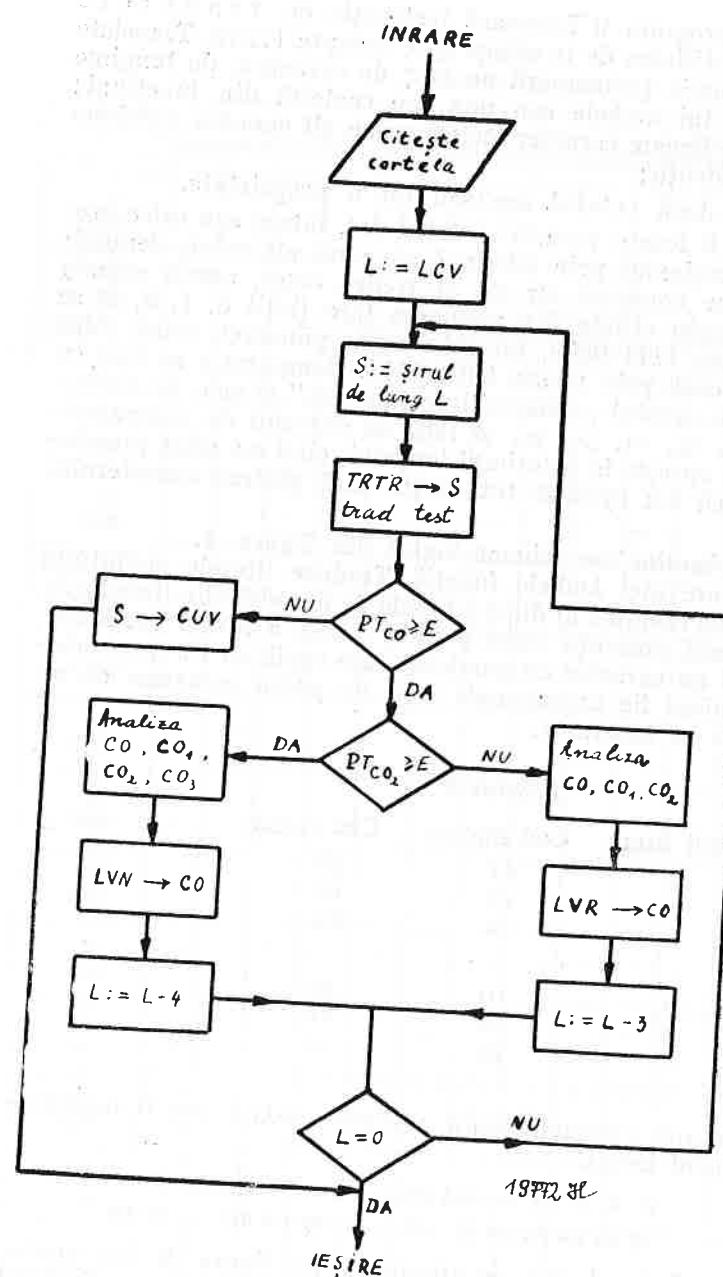


Figura 1. Rutina de codificare din EBCDIC în CGCLR

LCV = lungimea cuvintului
S = sirul ce se prelucrăază prin TRTR
CO = caracter de oprire
CO₁, CO₂ = caractere următoare
PT = ponderea tare
LVR = lit. văr. regulată
LVN = lit. var. neregulată
CUV = zona de depunere din memorie

LRD = lungimea radicalului
NAL = nr. alomorfelor radicalului
PT_{co} = ponderea tare a caracterului de oprire
AL_i = alomorful radicalului de rang *i*

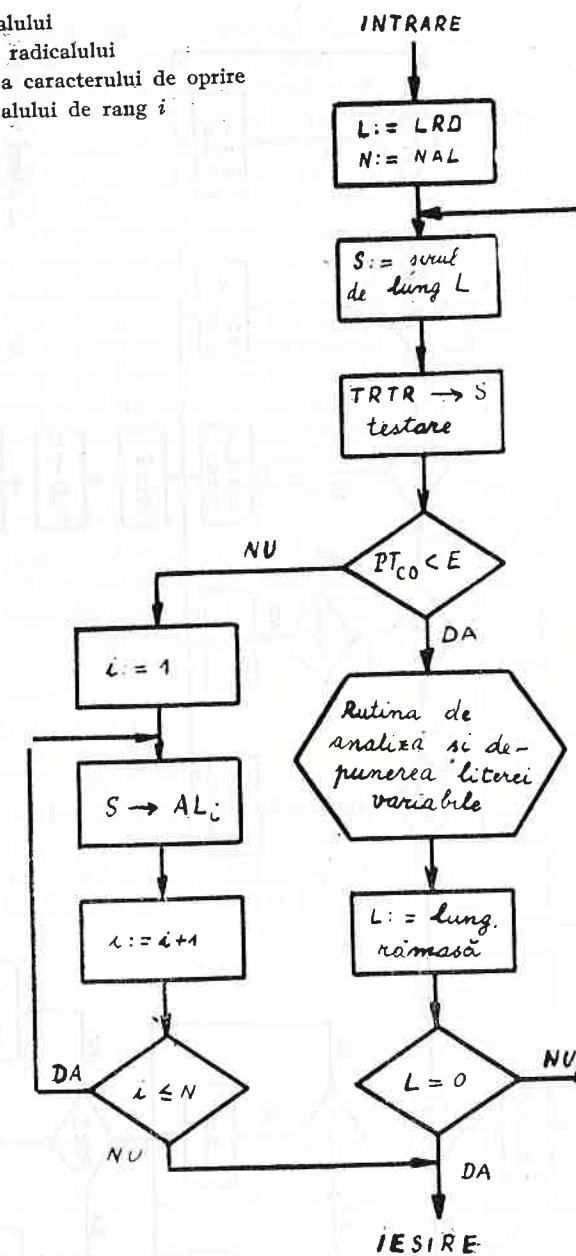


Figura 2. Rutina de formare a alomorfelor radicalului.

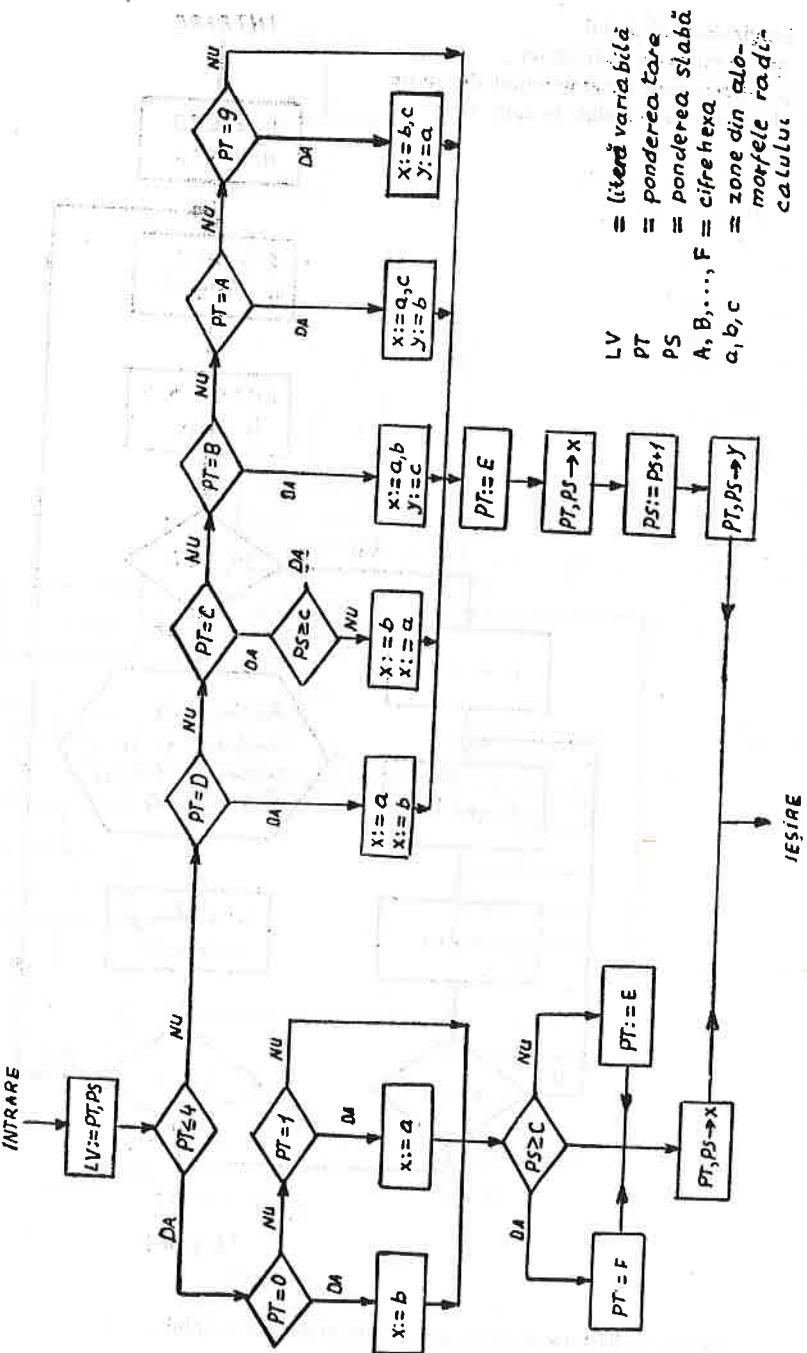


Figura 3. Rutina de analiză și depunere literelor variabile.

Instrucția TRTR testează (fără traducere) dacă șirul codificat în CGCLR conține litere variabile, cu pondere tare mai mică decât E. În caz afirmativ se analizează caracterul de oprire, adică litera variabilă aflată, dacă e regulată, dacă comută cu zero, etc., testând pondere tare și ponderea slabă. *Literele regulate ce nu comută cu zero permit o simplă formare a realizărilor*. Acestea sunt întotdeauna două, ambele din coloana E și una pe aceeași linie, alta pe linia următoare (tabelul 1). Deci cele două realizări vor fi E, PS și $E, PS + 1$ (PS este ponderea slabă a literei variabile). *Literele regulate ce comută cu zero au realizarea nevidă* pe aceeași linie, în coloana E (sau F, dacă sunt în zona punctată). *Realizarea vidă* impune mutarea spre stînga cu un pas a lanțului următor. Cele două realizări ale literei variabile se depun în alomorfele AL_i conform configurației, care de fapt servește la calculul adreselor din interiorul alomorfelor.

În această etapă cuvântul **ȚARA** (EC DO F5 E1) prezintă caracterul testat DO cu ponderea tare $D < E$. Fiind un caracter cu variație regulată DO devine întîi E_0 iar apoi devine E_1 . Cele două alomorfe sunt:

$$\begin{aligned} EC\ E_0\ F5 &= \ddot{T}\ A\ R \\ EC\ E_1\ F5 &= \ddot{T}\ \ddot{A}\ R \end{aligned}$$

Etapa II_B are schema logică în figura 4.

Printr-un ciclu în ciclu se depune fiecare alomorf AL_i în zonele P_{ik}^g din paradigmă — ca în schema (*) — potrivit cu partea de vorbire PV (coloana 33) cu gruparea GRP (col. 31, 32) care de fapt

- PV = partea de vorbire
- GRP = gruparea alomorfelor radicalului
- CLF = Clasa flexionară
- AL_i = alomorful de rang i
- P_{ik}^g = zona din paradigmă conform grupării GRP asociată radicalului AL_i
- FL_{j,c} = flectivul din clasa CLF asociat zonei P_j din paradigmă
- N = numărul alomorfelor radicalului
- NEP = numărul total al elementelor paradigmăi

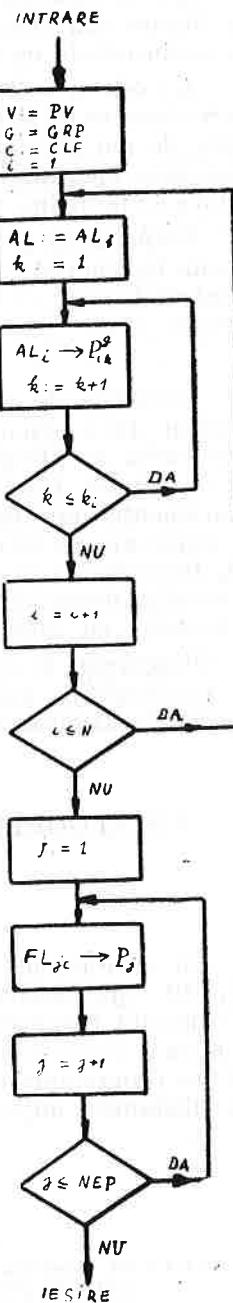


Fig. 4. — Rutina de atașare a flectivelor.

servește la calculul adreselor *zonenelor* p_{ik} din interiorul paradigmăi. Apoi la fiecare din cele N alomorfe ale radicalului AI_i se adaugă flectivele $FL_{j,i}$ în conformitate cu clasa flexionară CLF .

La cele $N = 2$ alomorfe ale radicalului cuvântului TARĂ, care au gruparea $GRP = 02$ se atașează seria de $2 + 4$ flective din declinarea numerotată de noi 13, care sunt -A -Ă și -I, -II, -ILE -ILOR. Ele se depun în cele șase zone ale paradigmăi cu ajutorul instrucției *transfer de caractere MVSR* (Move byte string right).

Etapa IIIc prin instrucția TRTR traduce fără testare din codul CGCLR în EBCDIC lanțurile de caractere ce sunt depuse la paradigmă după o tabelă inversă cu tabela 3 și apoi tipărește toată paradigmă cu cele NEP cuvinte aflate.

§ 6.

Cuvintele de pe cartele pot alcătui un *fisier pe disc* ce-l vom denumi *Dicționarul morfolologic* al limbii române despre care este vorba în [9] și [10]. Dicționarul permite *sinteză* formelor flexionare ale cuvintelor. Dotat cu un algoritm de analiză el permite și recunoașterea apartenenței unui cuvînt dat la o unitate lexicală. De exemplu, recunoaște că POPOARELOR aparține cuvântului POPOR și este CM6, pl., [gen., dat.,] art. Precizăm că proprietățile codului CGCLR ne permit o asemenea analiză în condiții destul de simple. Orice altă operație de sinteză a frazelor limbii, a textelor cu ajutorul ordinatorului utilizează acest dicționar.

Programul în ASSIRIS al algoritmului descris mai sus se află în curs de verificare pe calculatorul IRIS 50 de la Centrul Teritorial de Calcul Electronic Timișoara.

UN ALGORITHME POUR LA GÉNÉRATION DES MOTS DE LA LANGUE ROUMAINE

RÉSUMÉ

En utilisant les lettres à valeurs variables, conçues par Gr. C. MOISIL [9], [10], qui décrivent les alternances phonétiques des radicaux des mots on construit l'alphabet généralisé de la langue roumaine écrite. On a codé les mots à l'aide de cet alphabet en CGCLR (tableau nr. 1) et le décodage par un programme en ASSIRIS permet d'obtenir les formes flexionales des substantifs, adjektifs et verbes avec l'ordinateur IRIS 50 (FELIX C-256)

BIBLIOGRAPHIE

- [1] Minerva Bocșa, *Codage de l'alphabet généralisé du roumain écrit pour l'ordinateur IRIS 50 (FELIX C-256). Le substantif.* — Cahiers de linguistique théorique et appliquée, 2, 1973.

- [2] Minerva Bocșa, *Codage de l'alphabet généralisé du roumain écrit pour l'ordinateur IRIS 50 (FELIX C-256). L'adjectif.* — Cahiers de linguistique théorique et appliquée, 1, 1974.
- [3] Documentation de l'ordinateur IRIS 50.
- [4] Documentația calculatorului electronic FELIX C-256.
- [5] Gramatica limbii române, Editura Academiei Republicii Populare Române (București), 1963.
- [6] Valeria Guțu Romalo, *Morfologie structurală a limbii române, (Substantiv, adjetiv, verb).* Editura Academiei Republicii Socialiste România (București), 1968.
- [7] — *Indreptar ortografic, ortoepic și de punctuație,* Editura Academiei Republicii Populare Române (București), 1965.
- [8] Solomon Marcus, *Limbă și cod.* Studii și cercetări lingvistice, 1, 1966.
- [9] Grigore C. Moisil, *Probleme puse de traducerea automată. Conjugarea verbelor în limba română scrisă.* Studii și cercetări lingvistice, 1, 1960.
- [10] Grigore C. Moisil, *Problèmes posés par la traduction automatique. La déclinaison en roumain écrit.* Cahiers de linguistique théorique et appliquée, 1, 1962.
- [11] Brica Nistor Domokos, *Algoritm de traducere automată din limba engleză în limba română.* Editura Didactică și Pedagogică (București), 1966.
- [12] Paul Schveiger, Vera Drondoe, *À propos de la vérification de la sous-routine de déclinaison mécanique des substantifs,* Revue Roumaine de Linguistique, XII, 1, 1967.

Facultatea de Matematică-Mecanică
a Universității din Timișoara

Primit la 26.VI.1973.