

# Prefață

Născută în a doua jumătate a secolului trecut, concomitent în Europa și în SUA, lingvistica matematică are acum o întreagă istorie iar dezvoltarea ei este în plin progres. Școala românească a avut încă de la început o prezență solidă în domeniu (realizând chiar lucrări de pionierat), așa cum se poate vedea dintr-o lucrare de sinteză apărută în anul 1978 <sup>1</sup>: peste 500 de titluri de peste 120 de autori.

Lucrarea de față se referă la unele aspecte cantitative și formale ale limbajului natural, cu atenție specială acordată limbii române. Sunt câteva decenii de când nu am mai avut, la Facultatea de Matematică a Universității din București, o teză de doctorat cu această tematică, în ciuda faptului că s-au prezentat numeroase teze referitoare la modelarea matematică a limbajelor, domeniu în care România se manifestă puternic pe plan mondial.

Sunt două părți, bine distincte, chiar dacă în strânsă legătură în această lucrare: prima, caracterizată prin obiectul ei, silaba, ocupă capitolele al doilea (Aspecte ale silabei) și al treilea (Un model contextual-Marcus al silabei); a doua, caracterizată prin atenția îndreptată spre anumite probleme generale de similaritate, agregare și clasificare, dar care și ele sunt confruntate în final cu situații punctuale din studiul limbilor naturale, ocupă capitolele al patrulea (Despre similaritate), al cincilea (Metode multicriteriale de decizie) și al șaselea (Aplicații și rezultate experimentale). O bogată bibliografie de peste 120 de titluri (cele mai multe din ultimii 10 ani) folosită cu maturitate de către autor, încheie lucrarea.

Silabei i s-a refuzat statutul de unitate structurală a limbii, fapt care a contrastat cu acceptarea acestui statut pentru fonem și morfem. În consecință, nici modelarea matematică a silabelor nu a ajuns la complexitatea modelelor propuse pentru fonem și morfem. În deficit la înțelegerea calitativă, logică, a problemei silabei, cercetarea s-a orientat cu precădere spre aspectele cantitative, statistice. De exemplu, V. V. Ivanov atrăgea atenția asupra poziției mediane pe care silaba o are, între fonem și morfem: dacă numărul fonemelor este, în orice limbă, cuprins între  $10^1$  și  $10^2$  iar cel al morfemelor între  $10^3$  și  $10^4$ , numărul silabelor se află în intervalul dintre  $10^2$  și  $10^3$ . Există deci

---

<sup>1</sup>Marcus, S. Mathematical and computational linguistics and poetics. *Revue Roumaine de Linguistique*, XXIII, 559-588, 1978

o armonie combinatorială care urmează principiul minimului efort, aflat la baza multor legități cantitative care guvernează limbajul uman. Dar detaliile acestui fenomen se lasă cu greu identificate iar cercetarea din ultimele decenii încearcă să sporească lumina în aceste direcții.

Tocmai de aceea, Liviu P. Dinu întreprinde, în capitolul al doilea, un tur de orizont al ipotezelor și rezultatelor propuse de diverși autori. Nu cunoaștem să existe în literatura recentă o altă privire, la fel de bogată, asupra acestei chestiuni, dar, desigur, o atare sinteză rămâne totdeauna deficitară, atâta vreme cât nu se poate raporta la situația din fiecare limbă naturală. În ultima parte a capitolului al doilea, autorul propune un algoritm propriu de segmentare silabică a cuvintelor din limba română, după care este analizată una dintre puținele încercări, numai parțial reușite, de formalizare a silabei, încercare aparținând lui Theo Venneman (1978). Este astfel pregătit terenul pentru ceea ce se va întreprinde în capitolul al treilea, în care autorul prezintă propria sa încercare de modelare matematică a silabei lexicale (grafice).

Ceea ce autorul întreprinde în capitolul al treilea constituie o inițiativă dublă: mai întâi, propunerea, în studiul silabei, a folosirii unui instrument pe care nimeni nu s-a gândit până acum să-l asocieze cu o problematică de acest fel; în al doilea rând, perfecționarea instrumentului respectiv, cu consecințe posibile în domeniul teoriei gramaticilor formale. Gramaticile contextuale la care recurge Liviu P. Dinu sunt de natură predominant calitativă și utilizarea lor în studiul silabei confirmă faptul că problemele cantitative sunt organic legate de cele calitative. Cantitatea și structura se află într-o dialectică foarte fină, ele interacționează. Ideea ingenioasă a autorului pleacă de la analogia pe care o observă între tăietura silabică și procesul de generare în gramaticile contextuale; pe această bază, Liviu P. Dinu introduce o variantă nouă de gramatici contextuale totale și pe care le consideră în două variante. Cu ajutorul lor se definesc gramaticile silabice, care se plasează sub aspect generativ între gramaticile contextuale interne cu selecție și cele contextuale totale. Itinerarul propus de Liviu P. Dinu face jonțiunea cu așa numitele automate "go-through" introduse în 2001 de Radu Gramatovici (în cu totul altă ordine de idei), identifică astfel o clasă de automate eficiente pentru despărțirea în silabe și, pe parcurs, implică, în mod surprinzător, structuri de tipul șirurilor lui Fibonacci. Aceste structuri sunt introduse pentru a investiga cuvintele regulate din limba română, cuvinte care sunt construite pe un alfabet de două elemente, V (vo-

cală) și C (consoană); ele amintesc de problema pe care și-o punea A. Markov, în urmă cu aproximativ o sută de ani, când, căutând să modeleze alternanța consoanelor și vocalelor în "Evgeni Onegin" de Pușkin, a introdus noțiunea pe care mai târziu matematicienii au numit-o "lanț Markov". Și sub acest aspect, demersul lui Liviu P. Dinu este foarte atractiv. Dar miza sa este, în această privință, mai ambițioasă, prin apropierea pe care o face de teoria lui Levelt și Indefrey, din 2001, privind producerea vorbirii, teorie în care silabificarea constituie o etapă esențială. Plecând de la ipoteza unei funcționări mai degrabă paralele decât secvențiale a creierului în formarea frazelor, este propusă și o abordare paralelă a segmentării silabice prin introducerea gramaticilor de inserție cu derivare paralelă. Științele cognitive sunt aici direct vizate.

În capitolul al patrulea se pleacă de la tendința limbajului natural de a plasa informația cea mai semnificativă în prima parte a mesajului și se propune, într-o primă fază, o măsură de similaritate a unor clasamente diferite ale acelorași obiecte. Se consideră și cazul în care obiectele considerate în cele două clasamente sunt numai parțial aceleași. Sunt valorificate idei dintre cele mai variate, de la clasică distanță a lui Hamming din teoria codurilor până la măsurile de similaritate folosite recent în studiul genomului uman. Autorul introduce ceea ce numește *distanța rang*, permițând compararea unor clasamente de obiecte și lungimi diferite. Extensiunea acesteia la cuvinte permite definirea unei măsuri de similaritate a arborilor. Sunt identificate viitoare extinderi ale acestei măsuri, dintre care remarcăm extinderile posibile la calculul cu cuvinte al lui Zadeh, rezultatele putând fi astfel adaptate la operarea cu percepții.

Agregarea și clasificarea clasamentelor fac obiectul capitolului al cincilea. În raport cu distincția care se operează în literatură între metodele (de clasificare și agregare) supervizate (care necesită o mulțime de antrenament) și cele nesupervizate, autorul prezintă două metode supervizate de clasificare și propune propria metodă nesupervizată de clasificare. Toate aceste metode se bazează pe interferența deciziilor mai multor clasificatori, autorul demonstrând astfel că este conectat la tendința ultimilor 10-15 ani din literatura științifică de specialitate. Problema determinării autorului unui text controversat se situează printre aplicațiile posibile ale metodelor de clasificare multicriterială. Metodele de agregare și clasificare introduse de autor se prevalează de distanța rang introdusă în capitolul anterior și sunt comparate cu alte metode. Mai este

propus un algoritm de agregare polinomial și unul de clasificare, corespunzător metodei autorului. Urmează o analiză interesantă a modului în care metoda de agregare satisface una sau alta din condițiile de raționalitate propuse în această privință. Reperetele sunt constituite aici de faimoasa teoremă de imposibilitate a agregării a lui Arrow și de cele 17 condiții de raționalitate formulate de Gheorghe Păun. Sunt formulate, în încheiere, câteva probleme deschise.

Din ultimul capitol, al șaselea, dedicat unor aplicații, se detașează ca interes studiul comparativ al limbilor romanice, din punctul de vedere al celor mai frecvente silabe în fiecare dintre ele. Cuvintele din vocabularele reprezentative ale acestor limbi au fost descompuse în silabe, acestea din urmă fiind ordonate după frecvență, după care ierarhiile obținute au fost studiate din punctul de vedere al distanței rang dintre ele. O serie de tabele și grafice prezintă rezultatele obținute. Se constată că dintre distanțele rang ale limbii române față de celelalte limbi romanice cea mai mică este distanța față de limba portugheză. Pe de altă parte, se constată că, din punctul de vedere al structurii silabice, distanța limbii române față de una sau alta dintre limbile romanice occidentale este mai mare decât distanța dintre două limbi romanice occidentale, oricare ar fi acestea. Aceste rezultate concordă cu așteptările care provin din alte considerente, de ordin istoric și structural, dar aduc o serie de precizări de finețe.

Acesta este, în esență, conținutul lucrării propuse de dl. Liviu P. Dinu; lucrarea demonstrează capacitatea autorului de a întreprinde o cercetare conformă cu cerințele actuale ale științei, deci bazată pe o informație bogată și la zi și răspunzând unor probleme semnificative, în imediata continuare a cercetărilor anterioare ale altor autori sau ale sale proprii. După cum s-a putut vedea, autorul are nu numai cunoștințe ci și aptitudinea de a lua inițiative fericite, de a imagina procedee ingenioase.

Sperăm ca această carte să stimuleze interesul studenților și al comunității științifice cu preocupări în domeniul lingvisticii matematice și computaționale.