

ACCADEMIA NAZIONALE DEI LINCEI

ANNO CCCLXXII - 1975

E. S. S. 61  
DOPPIO

CONTRIBUTI DEL  
CENTRO LINCEO INTERDISCIPLINARE  
DI SCIENZE MATEMATICHE E LORO APPLICAZIONI  
N. 13

---

COLLOQUIO SUL TEMA:

LE TECNICHE DI CLASSIFICAZIONE  
E LORO APPLICAZIONE LINGUISTICA

(Firenze, 13 dicembre 1972)

(*ESTRATTO*)



ROMA  
ACCADEMIA NAZIONALE DEI LINCEI

1975

ANTONIO ZAMPOLLI

## L'ELABORAZIONE ELETTRONICA DEI DATI LINGUISTICI: STATO DELLE RICERCHE E PROSPETTIVE

Lo scopo di questa comunicazione non è quello di passare ordinatamente e sistematicamente in rassegna le diverse applicazioni dei calcolatori alla elaborazione dei dati linguistici. Non solo tale rassegna esigerebbe una classificazione di queste attività che per diversi motivi non sembra utile<sup>(1)</sup> ma andrebbe ben oltre il tempo concessoci. Preferisco invece discutere alcune ricerche, attraverso le quali è possibile, a mio avviso, illustrare non solo i risultati conseguiti negli ultimi anni, ma anche e soprattutto la natura dei problemi che oggi sono all'attenzione dei ricercatori e la situazione di stallo conseguente alla mancata soluzione di alcuni di essi. Le soluzioni di questi problemi costituiranno probabilmente le mete degli sforzi dei ricercatori nei prossimi anni e determineranno le linee di ricerca e le prospettive di sviluppo.

Mi sono basato soprattutto sulle mie esperienze di collaborazione ai numerosi progetti di spoglio condotti dai diversi utenti della Divisione Linguistica del CNUCE e sulle idee e sulle esperienze esposte dai rappresentanti delle mag-

(1) Una classificazione delle diverse applicazioni, dal momento che le tecniche per la elaborazione elettronica dei dati sono ormai impiegate praticamente in tutti i settori della ricerca linguistica, dovrebbe riflettere in realtà una classificazione delle ricerche linguistiche. Il panorama di queste ricerche oggi presenta, accanto a uno straordinario fiorire di attività, non poche contraddizioni, al punto che non sembra possibile tentarne una classificazione se non ponendosi nella prospettiva di una posizione teorica legata a una corrente o a una scuola ben determinata. E infatti gli autori che hanno fino ad oggi tentato una rassegna delle attività della linguistica computazionale (CL) hanno cercato invano, a mio giudizio, di sottrarsi a questo vincolo. Si può mostrare che anche coloro i quali – provenendo da discipline non linguistiche come la scienza dell'informazione, la logica, la matematica, l'ingegneria, la cibernetica – hanno adottato un criterio di classificazione dichiaratamente non linguistico, di fatti si sono mossi nel quadro di una teoria linguistica determinata. È il caso per esempio delle classificazioni operate in base alla « profondità » del trattamento dei dati e della complessità con cui opererebbero i diversi algoritmi, nelle quali la profondità viene misurata per confronto con i modelli generativisti di stampo semanticista, e il grado di utilità, di adeguatezza e di scientificità delle elaborazioni viene ritenuto direttamente proporzionale alla complessità degli algoritmi e all'impiego di strumenti matematici e formalizzanti. In realtà, molti cultori della CL, provenienti dalle discipline tecnico-scientifiche citate, hanno recato un grave danno alla credibilità scientifica della CL, perseguendo e valutando le applicazioni non per la loro utilità nei confronti della linguistica, delle discipline letterarie, della filologia, ecc., ma per la complessità dell'algoritmo e per il grado di formalizzazione che esso richiede. Per una bibliografia delle classificazioni proposte si veda Zampolli (1969 e 1970).

in altri Centri (Besançon, Leida, Cambridge, ecc.) che iniziavano la produzione di indici con sistemi UR a schede perforate (3).

Restava tuttavia praticamente senza risposta, nonostante alcune procedure costosissime e poco felici di ripiego, l'esigenza di produrre concordanze di struttura adeguata.

### 1.1.2. *Introduzione dei Calcolatori.*

Verso la fine degli anni '50 i calcolatori elettronici, già diffusi in campo commerciale e industriale da più di 10 anni in America e in Europa, affiancarono le macchine UR nelle elaborazioni per lo spoglio di testi.

L'adozione della nuova tecnologia non avvenne senza perplessità da parte dei lessicografi, e fu il principale argomento di discussione nei convegni di Tübingen (1960) e Besançon (1961). La scheda perforata, che può essere letta, maneggiata, e seguita costantemente di operazione in operazione attraverso la macchina UR, in fondo continuava ad essere vista come una normale scheda lessicale manuale, strumento tradizionale di generazioni di ricercatori.

L'uso dei calcolatori si diffuse in seguito, più rapidamente di quanto era stato previsto nei due congressi citati, ma gli sforzi dei lessicografi furono rivolti, in ultima analisi, a far eseguire al calcolatore le operazioni più pesanti della procedura tradizionale, senza modificarla sostanzialmente.

Il Congresso di Praga (1966) consacrò, per così dire, con accenti trionfalistici, l'avvento delle « machines dans la linguistique ». Il panorama dei centri attivi un pò ovunque in Europa, la varietà delle applicazioni e le dimensioni dei progetti suscitavano un grande entusiasmo. Era disponibile una ricca gamma di procedure e di programmi per produrre con il calcolatore indici di frequenza, rimari, indici inversi, concordanze, schede contestualizzate, ecc.

I principali progetti di dizionari storici avevano adottato o stavano adottando gli spogli elettronici, e in tutto il mondo si moltiplicavano gli spogli di testi, così che era pienamente giustificato parlare di *The death of the hand made concordance* (Raben, 1969) (4).

(3) Per questo periodo si possono consultare i primi numeri dei *Cahiers de Lexicologie* e De Tollenaere (1963).

(4) Come indicazioni bibliografiche generali per questo periodo si possono citare gli Atti di alcuni Congressi « Colloque International sur la Mécanisation des Recherches Lexicologiques », Besançon, 1961, in *Cahiers de Lexicologie*, 3 (1962); « Literary Data Processing Conference », Yorktown Heights (New York) (1964); « Le applicazioni dei Calcolatori Elettronici alle Scienze Morali e alla letteratura », Milano, 1962; « Les Machines dans la Linguistique », Praga, 1966; « L'Automazione Elettronica e le sue implicazioni scientifiche, tecniche e sociali », Accademia dei Lincei, 1968; « L'elaboration électronique en Lexicologie et en Lexicographie », Pisa, 1970; « Lexicon electronicum latinum », Pisa, 1968) e alcune riviste sorte nel frattempo per rispondere alle esigenze di diffondere informazioni in questo settore in così rapida espansione (oltre ai già citati *Cahiers de Lexicologie* del « Centre d'Etude du Vocabulaire Français » di Besançon, la *Revue* dell'« Organisation Internationale pour l'Etude des Langues Anciennes par Ordinateur » di Liegi, *The Finite String* della « Association for Computational Linguistics » degli USA, *Calculi* della « American Philological Association », *Computers and the Humanities* di New York, ITL dell'« Instituut voor Togepaste Linguistiek » di Lovanio, ecc.).

1.1.2.1. *Alcuni esempi.*

Mi sembra opportuno riferire qui rapidamente alcuni lavori a titolo di esempio.

I progetti sono così numerosi che, nonostante le inchieste condotte per molte lingue, siamo ben lontani da un censimento soddisfacente. Poiché i progetti in corso presso il CNUCE coprono ampiamente le diverse zone di questo settore, mi sembra opportuno scegliere gli esempi tra di essi.

*Archivi lessicali per la compilazione di grandi dizionari storici di una lingua.*

Il progetto più significativo è senz'altro la costituzione degli archivi per il *Tesoro della Lingua Italiana* delle origini e per il *Vocabolario Storico della Lingua Italiana*, ad opera dell'Accademia della Crusca<sup>(5)</sup>. Sono stati elaborati elettronicamente 30 milioni di occorrenze, il 40% delle quali sarà rappresentato nell'archivio. Negli ultimi 5 anni è stata pressoché completata la preparazione dell'input, ed ora le risorse umane e tecniche sono dedicate principalmente alla lemmatizzazione e alle scelte degli esempi. Le operazioni lessicologiche hanno raggiunto un notevole livello di formalizzazione, e formano una procedura che funge da modello per le elaborazioni degli altri progetti di spoglio, italiani e europei.

In questi lavori che affondano le loro radici nella ricca tradizione lessicologica e lessicografica europea, si riflette la tendenza moderna a creare grandi archivi ed inventari, che è stata determinata dalla diffusione dei calcolatori e dallo sviluppo della linguistica descrittiva. «The idea of creating lexicological archives is slowly making headway in each country. Slowly and cautiously without speaking of the naive myth of exhaustibility, these positive principles of cumulative work seem to be gaining the acceptance that they deserve (...) New methods of inventory-making have profoundly changed the conditions and the perspectives of lexicological text analysis» (B. Quemada, 1973, p. 425).

Gli esperimenti intrapresi su larga scala alcuni anni fa da centri nazionali, per es. a Besançon, Firenze, Nancy, Leida, si sono trasformati sotto un certo punto di vista in imprese di routine nelle quali lavorano grandi équipes di lessicografi.

Ma sembra chiaro che sta iniziando una nuova fase sperimentale: la quantità dei materiali raccolti — decine di milioni di citazioni — richiede che

(5) Cito come esempi di archivi lessicali di lingue con alfabeto non latino, preparati in Italia:

— il Progetto di un Dizionario Storico della Lingua Rumena del XVI secolo, i cui lavori, in collaborazione tra il CNUCE di Pisa e Bucarest, sono diretti da F. Dimitrescu (1973);

— il Lessico dell'ittita cuneiforme, diretto dal Prof. P. Meriggi (1973) di Pavia;

— l'archivio dell'assiro babilonese in collaborazione tra L'UCLA, il CNR e il CNUCE (Buccellati 1974).

il calcolatore venga adoperato non solo nell'assemblaggio degli archivi, ma anche come aiuto nella loro gestione. Avremo modo di discutere questa situazione più avanti.

*Archivi lessicali per la redazione di Dizionari Storici relativi a discipline particolari.*

A titolo di esempio cito il progetto per il « Vocabolario Storico della Lingua Giuridica » (Ciampi, 1973).

Questi progetti sono caratterizzati dalla esigenza di effettuare una forte selezione degli esempi, perché il lessico che deve essere messo in evidenza e accolto nell'archivio è un lessico nettamente specializzato: cioè il rapporto percentuale tra occorrenze in *input* e occorrenze che diventano *items* dell'archivio è molto basso. Perciò questi lavori hanno l'aspetto di esperienze interessanti sia metodologicamente sia economicamente, anche in relazione alla attuale richiesta di dizionari terminologici che proviene da numerosi settori scientifici e tecnici.

*Lessici, indici e concordanze dell'opera omnia di un autore.*

Merita rilevare che particolarmente numerosi sono i progetti relativi all'opera omnia di filosofi e pensatori. L'indice del lessico è considerato non solo come la chiave che permette di ricercare agevolmente i passaggi ove sono definite o esplicitate le idee, ma anche come uno strumento per chiarire e verificare l'interpretazione. Si pensi per esempio all'*Index Thomisticus* già citato, che dovrebbe iniziare entro l'anno la pubblicazione del primo dei 60 volumi previsti, e all'archivio dell'opera omnia di A. Rosmini, testè completato.

*Spogli di opere diverse.*

L'inventario più completo delle opere sottoposte a spoglio elettronico è probabilmente costituito dalle Sezioni « Directory of Scholars Active » e « Verbal Materials in Machine-Readable Form » che appaiono periodicamente in « Computers and the Humanities ». In Italia sono in corso di spoglio testi di diverse lingue (ventotto). Si tratta per lo più di studi stilistici, effettuati spesso per tesi di laurea su base lessicale e statistica, aventi come strumenti preferiti concordanze e indici. Le metodologie sono ben note, ma può essere interessante segnalare qualche esperienza di concordanze parzialmente diverse dalle forme tradizionali.

Le concordanze che chiamiamo *contrastive* si sono rivelate utili per confrontare un testo di base con un insieme di testi satelliti (per esempio un originale con le sue traduzioni in altre lingue, o con le sue successive redazioni). A ogni forma del testo originale si attribuiscono come contesti non solo il normale contesto, ma anche le unità contestuali corrispondenti nei testi satelliti. Analogamente, ogni forma del testo satellite riceve come contesto sia il proprio contesto normale, sia i contesti corrispondenti negli altri testi satelliti e nel

testo base <sup>(6)</sup>. Questo metodo ha permesso utili considerazioni critiche sul meccanismo della traduzione di opere letterarie. Le concordanze *sintattiche* hanno come esponente un classificatore o un insieme di classificatori dei tipi sintattici. Il contesto è costituito dalla proposizione che appartiene al tipo sintattico in esponente, accompagnata dalle proposizioni ad essa collegate da rapporti immediati di coordinazione, reggenza e subordinazione (cfr. De Mauro e Policarpi, 1971).

Le concordanze che chiamiamo *tematiche* ci sono state richieste per l'analisi delle strutture di testi poetici di tradizione orale. In queste concordanze l'esponente è costituito dalla sigla di un *tema* o di un *motivo*, e il contesto è l'intero brano corrispondente. Ovviamente sono prodotte anche tutte le possibili liste di correlazione tra lessico, temi, motivi, ecc. Il lavoro più avanzato in questo settore è l'analisi strutturale dei componimenti lirici di Pindaro e di Bacchilide.

Un particolare approfondimento metodologico richiedono gli spogli di testi condotti con interesse prevalentemente filologico. Il progetto di indici e concordanze di d'Arco Silvio Avalle (1973) sulla lingua poetica in Italia avanti la fine del XIII secolo è particolarmente interessante, perché presenta un ciclo completo di operazioni che, in interazione con il calcolatore, portano dalla trascrizione dei manoscritti, anteriori alla fine del XIII secolo, alla pubblicazione in fotocomposizione del testo riprodotto diplomaticamente e delle sue diverse analisi (indici, concordanze, ecc.).

Questo progetto è a mio avviso destinato a diventare un modello nel suo genere. Con una sola trascrizione del testo, e dunque con un solo controllo, si ottiene la stampa sia del testo critico sia degli indici e delle concordanze. Le correzioni eseguite con il metodo interattivo per mezzo del calcolatore sono più economiche e sicure delle correzioni di bozze. È dimostrato dalle esperienze di più ricercatori che gli indici e le concordanze sono strumenti molto efficaci per scoprire eventuali errori, sviste, incoerenze dell'editore di un testo. Le concordanze e gli indici che si producono ogni giorno al CNUCE mettono in evidenza errori e lapsus in edizioni critiche ritenute giustamente eccellenti e sicure. È dunque un grosso vantaggio poter pubblicare l'edizione di un manoscritto dopo averlo potuto analizzare con il calcolatore.

Un altro problema tipico è quelli del trattamento delle varianti. Oltre ad alcuni tentativi di automatizzare le fasi della *collatio* e della *recensio codicum*, che hanno condotto a risultati positivi per quanto concerne lo studio dei relativi algoritmi ma non sembrano ancora destinati ad applicazioni pratiche, va segnalato lo sforzo che alcuni filologi italiani stanno compiendo per definire dal punto di vista metodologico il trattamento delle varianti di un testo del quale si pubblicano indici, concordanze, ed elaborazioni statistiche. Si sono sperimentati metodi diversi. La soluzione più semplice prevede unica-

(6) Si veda, per esempio, la tesi di laurea di A. Penazzo (Parma, 1972) sulle diverse traduzioni italiane di alcuni brani del Faust di Goethe.

mente l'apposizione di un contrassegno alle parole per avvertire che hanno delle varianti che non vengono specificate. Soluzioni più complesse prevedono che tutte le informazioni dell'apparato critico vengano inserite nelle elaborazioni. In questo modo, le singole varianti e i rapporti tra le singole varianti sono presentati esplicitamente e perspicuamente nelle concordanze, negli indici, nelle statistiche. Il progetto più avanzato è lo spoglio delle tre redazioni dell'*Orlando Furioso* di L. Ariosto diretto da C. Segre <sup>(7)</sup>.

### *Banche di dati.*

L'archivio dell'*Atlante Linguistico Italiano* comprende circa 9 milioni di parole in trascrizione fonetica, raccolte con inchieste condotte in 1000 località italiane, in base a un questionario di circa 700 domande (cfr. Grassi 1973). Sono già state pubblicate alcune analisi di tale questionario (Ali). È in corso la trasposizione in 'machine readable-form' delle risposte dialettali. L'adozione delle tecniche elettroniche ha indotto a capovolgere l'ordine tradizionale delle operazioni, nel quale la messa in carta dei materiali precedeva gli indici e le analisi. Il piano attuale prevede invece la pubblicazione dello *Atlante* solo dopo la pubblicazione di indici regionali, di dizionari per località, di carte relative a fenomeni linguistici particolari, ecc. Si pensa già a un progetto di estensione dell'archivio con nuove inchieste, strutturate in una banca dinamica delle conoscenze dialettali. Si deve tener presente il rinnovamento dell'interesse per la dialettologia in Italia, motivato, tra l'altro, dall'urgenza di fissare, prima che vada perduto, il patrimonio culturale regionale, e dalla posizione privilegiata degli studi dialettali e regionali nel quadro delle ricerche sociolinguistiche, attualmente molto sentite nel contesto di mobilità sociale e geografica di gruppi e di individui.

#### 1.1.2.2. *Procedure di spoglio.*

Esaminiamo a grandi linee lo schema di uno spoglio lessicale, così come viene condotto di solito presso la maggior parte dei Centri Specializzati. Esiste naturalmente una grande varietà di procedure, ma esse hanno in comune alcune funzioni, ordinate, di base, alle quali qui mi riferisco.

#### Fase 1. — *Preparazione dell'input.*

Il testo viene registrato su un supporto operabile dal calcolatore. Le operazioni fondamentali sono le seguenti. Il testo viene ricopiato, da un operatore, di solito su schede meccanografiche. Le schede, o il loro equivalente

(7) È stato messo a punto un sistema che permette di comunicare al calcolatore il testo delle due prime redazioni semplicemente segnalando le differenze tra queste e la terza edizione. L'algoritmo di elaborazione decide in base a regole formali, che considerano la natura delle varianti, quali varianti porre in esponente, quali riportare nel contesto, nonché il loro conteggio nel calcolo delle frequenze.

meccanografico, vengono listate, cioè viene stampato dal calcolatore il testo così come è stato copiato dall'operatore. Questa lista viene letta e confrontata, di solito dal ricercatore, con il testo originale, per ricercare e segnare gli inevitabili errori di ricopiatura. Questi errori vengono corretti per mezzo di schede. Il testo così corretto viene registrato in forma definitiva su un nastro magnetico o su di un supporto equivalente.

Fase 2. — *Elaborazioni relative alle forme grafiche.*

Al calcolatore viene di solito richiesto di elaborare una serie di indici e concordanze nei quali l'unità di elaborazione è la *forma grafica*. In questa fase cioè la parola è definita, per il calcolatore, come una sequenza di caratteri alfabetici (lettere e segni diacritici: accenti vari, dieresi, apostrofo, ecc.) tra due spazi. Il calcolatore considera due sequenze identiche di caratteri alfabetici come due parole uguali, e cioè come due occorrenze di una stessa forma grafica. Perciò tutte le sequenze *faccia* del testo sono considerate come occorrenze di una stessa forma, anche se hanno significati diversi nei diversi luoghi del testo. Per la stessa ragione saranno tenute distinte due sequenze come *amavo* e *amerò*. Gli indici più comunemente elaborati sono l'index locorum (elenco delle forme grafiche ciascuna seguita dalle indicazioni di tutti i luoghi del testo ove appare), gli indici alfabetici diretto e inverso delle forme e delle relative frequenze, l'indice delle forme in ordine di frequenza decrescente, diversi tipi di rimario, le concordanze delle forme, e, raramente, le schede lessicali.

Fase 3. — *Lemmatizzazione.*

Sulle liste delle concordanze per forma, il ricercatore analizza linguisticamente le varie forme. Normalmente questa analisi si limita alla *lemmatizzazione*; il ricercatore, cioè, scrive accanto a ciascuna forma il lemma cui la forma deve essere assegnata. Se la forma è omografa, il lemmatizzatore indica il lemma per ciascuna delle sue occorrenze, esaminando il contesto di ciascuna occorrenza. Un operatore perfora poi queste indicazioni su schede per mezzo delle quali il calcolatore aggiunge il lemma a ciascuna parola del testo registrato su nastro magnetico. Di solito, come si è detto, l'analisi non arriva a livelli più approfonditi. Per esempio non viene operata a livello morfologico, per distinguere forme di uno stesso lemma omografe sul piano morfologico (del tipo *dica*), nè a livello sintattico, per distinguere le diverse 'costruzioni' o 'strutture' nelle quali la parola è inserita, nè a livello sematico, per distinguere le diverse accezioni di una parola polisemica. Spesso, nel corso della lemmatizzazione, viene effettuata una scelta dei materiali lessicali; cioè vengono scelti i lemmi e le occorrenze ritenuti 'interessanti': le altre parole vengono eliminate, e cioè ricevono un contrassegno che le esclude dalle elaborazioni successive (per esempio dalla produzione delle schede lessicali).

Fase 4. — *Elaborazione dei risultati dello spoglio.*

Si richiede al calcolatore l'elaborazione e la stampa degli stessi indici elencati alla fase 2, ma questa volta l'unità di elaborazione è il lemma.

Si pone a questo punto, di solito, il problema di mettere a disposizione di tutti gli studiosi questi risultati, pubblicando indici e concordanze dei lemmi. Farli ricomporre tipograficamente presenta l'inconveniente della correzione delle bozze, la cui gravità appare chiara se si considera che i risultati dello spoglio consistono in un numero di righe che è circa dieci volte quello del testo originale. Si preferisce riprodurli con metodi fotografici o xerografici, o ricorrere alla stampa per mezzo della fotocomposizione.

Le schede lessicali vengono prodotte dal calcolatore in due o tre serie complete, che vengono conservate in appositi archivi-schedario.

Tutta questa documentazione (indici, concordanze, e soprattutto le schede lessicali) è di solito alla base della redazione dei grandi dizionari storici, sincronici, ecc. È chiaro che quanto minore è il grado di analisi nel corso dello spoglio, tanto più oneroso e complesso sarà il compito del redattore. L'alternativa tra l'approfondire l'analisi e la scelta dei materiali in fase di spoglio, e il rimandarle alla fase di redazione, è oggi motivo di vivacissime discussioni tra lessicologi e lessicografi.

1.1.3. *La crisi attuale: Il problema centrale della lessicologia e della lessicografia assistite dai calcolatori.*

Verso il 1969-1970 cominciarono a levarsi le prime voci di insoddisfazione e di critica contro il generale ottimismo che, come abbiamo visto, caratterizzava, e caratterizza ancora oggi, l'ambiente degli utenti del calcolatore per gli spogli lessicali.

Ci si è resi conto che lo sviluppo delle applicazioni secondo le linee, per così dire, ormai 'classiche' degli anni '50 e '60 è giunto a un punto di saturazione. Se si continua secondo la metodologia corrente o, in altre parole, secondo le attuali regole del gioco, non esistono concrete prospettive di sviluppo. Le tecniche e le metodologie in uso permettono di fare eseguire al calcolatore alcune delle operazioni tradizionalmente riservate al lessicografo; essenzialmente, quelle connesse alla trascrizione dei contesti per ciascuna occorrenza del testo e al loro ordinamento alfabetico. Ma per quanto le macchine diventino sempre più veloci e i programmi più sofisticati, i lessicografi non possono profittarne in proporzione, perché le metodologie attuali producono già molti più dati di quanti una équipe redazionale di dimensioni ragionevoli possa elaborarne lavorando secondo le procedure attuali. Ci si trova di fronte a una alternativa abbastanza chiara. Si accettano i limiti delle applicazioni odierne, si continuano a registrare *corpora* sempre più vasti di testi, a produrre archivi di concordanze o di schede lessicografiche senza alcuna previa analisi e classificazione linguistica. Tali elaborazioni vengono rinviate a una successiva fase di redazione, senza preoccuparsi per il momento degli enormi pro-

blemi che saranno posti dalla quantità dei materiali lessicali archiviati. Oppure si rinnovano le metodologie attuali. Si applicano i ritrovati di alcune discipline correlate, quali la linguistica computazionale, la statistica linguistica, l'*information retrieval*, ecc. Si sfruttano adeguatamente le nuove possibilità che la tecnologia dei calcolatori propone nel suo rapidissimo sviluppo; in particolare le tecniche conversazionali che permettono di usare il calcolatore come uno strumento redazionale attivo, in stretta interazione con il lessicografo (8).

## 1.2. IMPIEGO DEI CALCOLATORI. NUOVE TENDENZE E PROSPETTIVE.

Esaminiamo ora le linee di ricerca secondo le quali si cercano di superare le difficoltà e i limiti dei metodi attuali.

### 1.2.1. *Preparazione dell'input.*

Le difficoltà sono sia di ordine tecnico (9) sia di ordine economico.

(8) In queste affermazioni ci troviamo d'accordo B. Quemada ed io nelle nostre comunicazioni al « Colloque sur l'élaboration électronique en lexicologie et lexicographie » da me organizzato a Pisa nell'estate del 1970 (Zampolli ed., 1973 a). Questa nostra presa di posizione fu confermata dalla « Tavola Rotonda sui grandi lessici storici » organizzata a Firenze presso l'Accademia della Crusca nella primavera del 1971. Essa propose di inserire nella 2<sup>a</sup> Scuola Estiva Internazionale di Pisa, accanto al già previsto indirizzo dedicato alle grammatiche formali e al trattamento automatico di strutture sintattiche e semantiche, un nuovo indirizzo per studiare in che misura e in che forma il calcolatore possa assistere il lessicologo e il lessicografo anche nella fase di analisi dei dati lessicali e di redazione, oltre che nella loro raccolta. Lo svolgimento contemporaneo dei due indirizzi ha permesso di osservare concretamente quanti temi e problemi in comune abbiano lessicologia e lessicografia da un lato, e linguistica computazionale e molte correnti teoriche attuali della linguistica dall'altro. È risultato molto chiaro che lo studio del lessico si prospetta oggi anche per queste discipline come un tema centrale. I linguisti si rendono conto che è indispensabile affrontare il compito di applicare, in estensione, i sistemi di classificazione e le tecniche di analisi e di descrizione a sottoinsiemi lessicali sempre più vasti, e sono portati a rivalutare l'opera di raccolta dei dati, portata avanti fino ad oggi quasi esclusivamente dai lessicografi, e le loro esperienze. I lessicografi, dal canto loro, prendono coscienza del fatto che l'interesse dei linguisti teorici e dei linguisti computazionali, convergendo sul lessico, ha prodotto strumenti teorici e tecnici che rispondono ad esigenze fondamentali del compito lessicografico. Questo punto di vista è stato chiaramente espresso anche da L. Venezki e W. P. Lehman alla « International Conference on English Lexicography » del 1972 a New York (cfr. Zampolli, 1973 c) ed ha trovato il consenso di quanti, linguisti e lessicografi, hanno partecipato al « Deuxième Colloque International de Linguistique et de Traduction » (Montréal 1972) cfr. Zampolli, 1973 d e al « Colloque sur l'analyse des corpus linguistiques: problèmes et méthodes de l'indexation maximale » (Strasburgo, 1973). Un'ulteriore conferma di questo atteggiamento è costituito dal fatto che le applicazioni dei calcolatori allo spoglio di grandi corpora è stato posto per la prima volta in modo esplicito tra i temi della « 5 TH International Conference on Computational Linguistics » di Pisa (1973).

(9) È noto infatti che le tastiere delle macchine abitualmente disponibili per ricopiare il testo su supporti leggibili dal calcolatore sono così povere di caratteri da rendere necessaria l'indirizzo e gli eventuali nuovi caratteri di correzione; immettere le correzioni nel calcolatore;

L'operazione di ricopiare il testo è di per sè lunga e molto costosa. Il ciclo di controllo e di correzione, che secondo la nostra esperienza deve essere ripetuto in media 3 volte, triplica il tempo e il costo. Per preparare 3 pagine di circa 40 righe ciascuna, occorrono di solito circa 3 ore così ripartite: perforazione 1 ora; verifica 50'; lettura per controllo 30'; perforazione ed esecuzione delle correzioni 20'; due cicli ulteriori di controllo e di correzione 20'. Non vanno poi dimenticati i 'tempi morti' per il passaggio dei materiali tra i diversi esecutori del ciclo (controllore, perforatore, calcolatore). I rimedi a questa situazione possono essere sia tecnici sia organizzativi.

L'evoluzione tecnologica sta per rendere obsolete le perforatrici di schede o di nastro: per una documentazione completa dei nuovi mezzi disponibili rinvio all'articolo di B. R. Schneider (1971), soprattutto per le possibilità offerte dai lettori ottici e dai terminali <sup>(10)</sup>.

una codificazione difficile e gravosa per rappresentare la ricca varietà di caratteri presenti nei testi. Le complicazioni della codificazione aumentano il numero degli errori di battitura. Il ricercatore, nel tentativo di semplificarla, rinuncia spesso a rappresentare grafemi che non sembrano immediatamente indispensabili per le elaborazioni che egli ha progettato. Accade poi non di rado che il testo così impoverito non sia sufficiente per altre ricerche o altri ricercatori, i quali debbono provvedere a riperforarlo. La scarsità di caratteri disponibili nelle normali stampanti, nuoce, naturalmente, anche alla rapidità e alla esattezza del *controllo a lettura*.

(10) Qualcuno ha detto che nel campo dei lettori ottici l'utopia è sempre dietro l'angolo. Dapprima si è pensato a un lettore ottico per un solo *font*; oggi esistono lettori ottici *multi-font*, e già si lavora a un lettore ottico *omnifont*. Questo non esistesse ancora, ma è stato dimostrato che con tecniche appropriate è possibile aggiungere con poca spesa nuovi *fonts* a quelli leggibili da un certo sistema. Se si deve spogliare un corpus stampato per intero con un solo *font*, è consigliabile calcolare se convenga o no sostenere le spese della programmazione necessaria per aggiungere questo nuovo *font*. Comunque è sempre possibile battere il testo con una macchina da scrivere dotata di un *font* già accettato da un determinato sistema ottico, e di farlo poi leggere da questo. Sembra accertato, da esperienze condotte, che rispetto alla perforazione il costo scenda ai 2/5 e che gli errori siano meno di 1/5. Il grado di accuratezza in perforazione influisce poi sull'andamento dei controlli successivi, perché gli errori che sfuggono al controllo sono in proporzione alla quantità di errori presenti. Inoltre la registrazione del testo prodotta da un lettore ottico può essere una immagine della pagina stampata. Ogni carattere, ogni corpo, e la loro collocazione nella pagina sono registrati automaticamente. Ciò significa che è possibile rilevare e contrassegnare automaticamente le categorie di informazione indicate dal tipo di formato, come, per esempio, le diverse sezioni di un articolo di dizionario.

Un interessante sistema è quello dei terminali CTR (Cathode Ray-Tube Terminal). Si tratta di un terminale che scrive su uno schermo TV (anziché o oltreché su un foglio di carta), e tutto ciò che appare sullo schermo può essere trasmesso al calcolatore. Ci sono molti tipi, alcuni dei quali prevedono uno schermo opaco per non affaticare la vista e possiedono un insieme di caratteri molto ampio. Alcuni tipi sono anche programmabili, nel senso che l'insieme dei caratteri è prodotto dal *software*, e può essere moltiplicato a piacere. Il teleschermo è uno strumento nel quale il testo è perfettamente elastico; la cancellazione è istantanea e gli errori possono essere localizzati corretti e verificati con una sola operazione, anziché con le 6 operazioni che, nella procedura normale, sono: stampare il testo con un numero di ordine (indirizzo) univoco per ogni parola o riga; annotare l'indirizzo dell'errore e scrivere accanto la correzione relativa; perforare l'ordine opportuno (cambia in, cancella, inserisci) per ogni errore, il suo

Ancora di più è possibile fare sul piano organizzativo, nazionale e internazionale. In numerosi paesi, alcune case editrici hanno adottato il sistema di stampare con il *type-setting* o con la *fotocomposizione* e così producono una grande quantità di testi su nastro magnetico o di carta. In Italia stiamo già lavorando per assicurare la disponibilità di questi materiali per le ricerche linguistiche.

L'obiettivo più importante è però quello di creare degli organismi per coordinare a livello nazionale gli sforzi dei vari Centri e dei singoli ricercatori. È ovvio, innanzitutto, che occorre evitare i doppioni. Molti sarebbero sorpresi di sapere quante diverse versioni meccanografiche esistono di uno stesso testo, per es. Omero, la Bibbia, la Divina Commedia, e quanti programmi diversi operino lo stesso tipo di spoglio e di elaborazione sullo stesso tipo di macchine.

È necessario che i testi vengano registrati secondo gli stessi criteri tecnici e scientifici, così da poter essere facilmente scambiati tra i diversi ricercatori, e da poter costituire una grande biblioteca elettronica standardizzata, elaborabile con gli stessi programmi fondamentali. L'adozione di uno schema unico di registrazione assicura anche che vengano riprodotte tutte le informazioni presenti in un testo a livello grafemico, e quindi ne garantisce la utilizzabilità per le ricerche diverse. L'adozione di un tale *standard* permette ovvie economie di tempi e costi, e risponde a precise esigenze di ordine scientifico: i risultati dello spoglio di un testo dovrebbero prestarsi alla comparazione con gli altri testi dello stesso autore, della stessa epoca, della stessa scuola, dello stesso genere letterario, ecc.

### 1.2.2. « *Utility programs* ».

La standardizzazione degli archivi di testi è condizione necessaria per la standardizzazione dei programmi.

Ci sono oggi per esempio numerosissimi programmi di contestualizza-

indirizzo e gli eventuali nuovi caratteri di correzione; immettere le correzioni nel calcolatore; stampare le righe corrette; controllare che le correzioni siano state apportate nel luogo e nel modo voluti. Con il teleschermo invece l'utente semplicemente muove con un tasto il cursore, un contrassegno che indica e determina il carattere sul quale si può operare in un dato momento. L'utente può, operando con una tastiera, cancellare, sostituire, inserire determinati caratteri, parole, o intere frasi. Se qualcosa viene cancellato, lo spazio vuoto viene riempito automaticamente dal testo che si muove verso sinistra; se qualcosa viene aggiunto, il testo si sposta proporzionalmente verso destra. Alcuni terminali hanno anche la capacità di *formatting* che semplifica la codificazione nelle applicazioni (per esempio quelle bibliografiche o nelle predisposizioni di lessici) in cui una stessa struttura si ripresenta ripetutamente. Così, per esempio, la struttura appare sullo schermo come un modulo che deve essere riempito: autore, titolo, editore, città, data. L'operatore riempie gli spazi appositi. Ogni elemento è così classificato automaticamente in virtù della sua posizione sullo schermo prima di entrare in memoria. Un altro vantaggio è dato dal fatto che nel caso della battitura del testo o delle correzioni il calcolatore può integrare l'operatore applicando un programma che esercita dei controlli di correttezza formale: per esempio verifica se sono rispettate certe regole di fonotassi o consulta un lessico automatico. (cfr. Szanser 1969).

zione, cioè programmi che per ogni parola del testo (o per alcune categorie di parole del testo) producono un *record* costituito da almeno tre elementi: la parola, il suo riferimento al testo, il suo contesto immediato. Evidentemente, le differenze tra i diversi programmi consistono nell'algoritmo che ritaglia il contesto. Dopo aver esaminato i diversi tipi di algoritmi in uso presso i centri specializzati nello spoglio di testi, la *Divisione Linguistica del CNUCE* ha messo a punto un programma di contestualizzazione al quale l'utente può chiedere, con poche *control-cards*, di eseguire l'uno o l'altro tipo di algoritmo <sup>(11)</sup>. Questo lavoro di generalizzazione è stato fatto per tutte le fasi dello spoglio lessicale, dello spoglio fonetico, di alcune elaborazioni statistiche. Le esperienze condotte su circa 2.000 testi per circa 50 milioni di parole in 28 lingue diverse, ci consentono di affermare che le nostre procedure sono utilizzabili, grazie a poche schede-controllo, per elaborare qualsiasi testo in qualsiasi lingua, o almeno tutti i testi che possono essere ricondotti a una scrittura alfabetica.

### 1.2.3. *Lemmatizzazione semiautomatica.*

L'insieme delle operazioni comunemente raggruppate sotto il termine *lemmatizzazione* richiede una serie di interventi umani che, rompendo il ritmo delle elaborazioni interamente automatiche dello spoglio, aumenta considerevolmente il tempo, il costo, e i rischi di errore.

Molti ricercatori, in particolare anglosassoni e statunitensi, decidono di non lemmatizzare affatto i testi, e si limitano a produrre degli indici, delle concordanze, o delle schede-contesto nei quali gli esponenti non sono unità definite secondo criteri linguistici, ma semplici *forme grafiche*, e cioè 'parole' come le riconosce abitualmente il calcolatore: sequenze di lettere tra due spazi o tra due separatori in genere. Indubbiamente questa semplificazione ha qualche vantaggio. La velocità del calcolatore nell'operare su simboli viene sfruttata completamente e si possono produrre rapidamente e a costi minori grandi quantità di spogli che, se diffusi, rendono innegabili servizi agli studiosi. Alle volte ci sono ragioni scientifiche che sconsigliano la lemmatizzazione, per esempio nel caso di testi in lingue o strati di lingua poco noti, nei quali molti lemmi non sarebbero attestati o addirittura neppure ricostruibili.

Tuttavia, molto spesso, in particolare negli spogli di grandi corpora per la redazione di ampi dizionari storici, una qualche analisi e classificazione dei materiali lessicali sembra indispensabile prima della conclusione degli spogli, per evitare che i redattori del dizionario restino sommersi dalla quantità di dati da scegliere e da ordinare. È inevitabile chiedersi se il calcolatore, per aiutare effettivamente il lessicografo, non debba affiancarlo anche, e soprattutto, nella fase di classificazione dei dati lessicali che il calcolatore raccoglie in proporzioni non commesurabili alle possibilità umane di elaborazione.

(11) Si veda a questo proposito Zampolli (1972).

Il cosiddetto dizionario di macchina, o lessico automatico (LA), fornisce una prima risposta a questa esigenza.

1.2.3.1. *Il lessico automatico.*

Per maggiore chiarezza riassumo rapidissimamente cosa si intende con i termini LA e consultazione di un LA. Mi riferirò costantemente, per brevità di esposizione, alla più semplice tra le diverse possibili strutture di un LA. In nota è descritta una organizzazione più complessa <sup>(12)</sup>. Nella sua forma più semplice un LA consiste in una serie di *forme grafiche* registrate su nastro, su disco, o su altro supporto leggibile dal calcolatore. Queste forme sono ordinate alfabeticamente. Ogni forma è accompagnata da una serie di informazioni linguistiche di natura diversa a seconda dei diversi impieghi cui il LA è destinato. Chiamiamo *analisi* o *funzione* di una forma l'insieme delle informazioni che la accompagnano: per esempio se il LA è compilato come ausilio agli spogli lessicali, per ogni forma saranno dati il lemma cui la forma deve essere assegnata e, spesso, la classificazione grammaticale e morfologica del lemma e della forma. Se il LA è compilato per tradurre automaticamente da una lingua all'altra, saranno date anche la forma corrispondente nella lingua di uscita, alcune indicazioni per l'analisi sintattica e semantica della proposizione, ecc. Se il LA serve per statistiche fonematiche, etimologiche, sociolinguistiche, saranno dati anche la trascrizione fonemica della forma, la sua etimologia, i suoi registri d'uso, ecc.

(12) Un metodo di ricerca più complesso, usato soprattutto come input al *parser*, è quello della segmentazione. Nel LA sono registrate non tutte le forme grafiche di un lemma, ma solo il suo tema (o i suoi diversi temi, se necessario); il termine *tema* non deve essere inteso nel senso tecnico della linguistica, ma nel senso operativo di sequenza di grafemi che resta invariabile nel corso di tutta (o parte) la flessione. Il tema è accompagnato da un codice che rinvia a una tabella ove sono enumerate tutte le desinenze possibili, nel sistema morfologico considerato, per la categoria alla quale il tema appartiene. Parole invariabili e forme irregolari sono registrate nel LA come termini indipendenti. Se la parola da lemmatizzare non coincide con uno di questi termini, il programma cerca di segmentarla in *tema + desinenza*, in modo che il tema sia presente nel LA e che la desinenza faccia parte di quelle indicate nella tabella associata a tale tema. Naturalmente quando il programma trova più di una segmentazione corretta, tratterà la forma come omografa. I sistemi di segmentazione variano da algoritmi interessanti da un punto di vista puramente informatico ad algoritmi che incorporano una teoria linguistica morfologica specifica, e funzionano come interpreti di regole morfologiche e morf fonologiche. Per una discussione di questi algoritmi, particolarmente interessanti quando trattano oltre ai fenomeni di flessione, elisione e troncamento, anche quelli della derivazione mediante affissi e della composizione, si veda Kay (1974); il sistema di consultazione per segmentazione è indubbiamente più interessante dal punto di vista algoritmico, ma non è detto che sia sempre più economico della ricerca per forme. Per confrontare correttamente il rendimento dei due metodi, occorre tener conto dell'intera procedura di spoglio nella quale la consultazione è inserita, e delle dimensioni del testo da lemmatizzare. La ricerca per forme e non per temi conviene quando la procedura prevede in ogni caso la alfabetizzazione del testo e il suo rendimento cresce proporzionalmente al numero complessivo delle occorrenze. Il LAI (Lessico Automatico Italiano) esiste nelle due versioni, e cioè come inventario di temi (150.000 circa) e di forme (1.000.000 circa).

A seconda degli scopi cui il LA è destinato può variare sensibilmente anche il numero delle forme che lo compongono. Per esempio se il LA serve per tradurre testi scientifici relativi a una disciplina specifica, esso contiene di solito solo le voci più frequenti nella lingua comune e i termini tecnici della disciplina in questione. Se invece il LA serve per statistiche sull'intero sistema lessicale, l'insieme delle voci che compongono il LA è molto più ricco. Il LA si propone in tal caso come un sottoinsieme rappresentativo dell'intero lessico di una lingua, e al limite vorrebbe coincidere con esso <sup>(13)</sup>.

È importantissima, dal punto di vista applicativo, la distinzione tra forme univoche e forme omografe. Diremo *univoca* una forma alla quale, nel LA, corrispondono due o più analisi distinte. Per esempio se in un dato LA la analisi delle forme è rappresentata solo dal lemma, una forma come l'italiano *dica* sarà, in tale LA, univoca: essa può appartenere infatti solo al lemma *dire*. Sarà invece omografa nello stesso LA, la forma *amo*, che può appartenere sia al sostantivo *amo* sia al verbo *amare*. Consideriamo invece un LA nel quale l'analisi assegnata a ciascuna forma comprenda, oltre al lemma, anche la classificazione morfologica della forma (genere, numero, grado per le forme nominali; modo, tempo; persona per i verbi, ecc.). In un tale LA sarà omografa anche la forma *dica*, cui corrisponderanno 4 distinte analisi: rispettivamente 3<sup>a</sup> persona singolare dell'imperativo, 1<sup>a</sup>, 2<sup>a</sup>, e 3<sup>a</sup> persona singolare del congiuntivo presente. Se invece le analisi di un LA comprendessero solo la trascrizione in alfabeto fonetico delle forme, in tale LA forme come *amo* e *dica* sarebbero univoche, perché vanno in ogni caso trascritte rispettivamente come /ámo/ e /díka/, e sarebbero omografe solo le forme non omofone in italiano, come *ancora* (/ánkora/ - /ankóra/) e *pèsca* (/pèska/- /péska/). Come esempio possiamo riferirci al Lessico Automatico Italiano (LAI) che stiamo predisponendo al CNUCE sotto gli auspici del Comitato 08 del CNR.

Abbiamo registrato su nastro magnetico circa 150.000 lemmi, ottenuti dall'unione delle nomenclature dei principali dizionari. Ad ogni lemma sono state assegnate informazioni riguardanti l'etimologia, la funzione grammati-

(13) Si pone il problema se il lessico di una lingua debba o no considerarsi finito. Spesso la risposta è diversa a seconda che si tratti di morfemi o di parole. Si veda per esempio Bloomfield (1933, cap. 6) Hjelmslev (1953, cap. 12 e 14) Spang-Hanssen (1967) Wagner (1967, p. 17) Ch. Muller (1968, p. 134). L. Guilbert (1967, pp. 116 e segg.) e Rey-Debove (1970, pp. 3 e segg.) hanno discusso il problema dal punto di vista generativista. Si cfr. anche Zampolli (1973 b, p. 165) Un LA può essere compilato con procedimenti diversi. Per una lingua il cui lessico è ben noto, si registrerà dapprima un lemmario ottenuto, di solito, dall'unione delle nomenclature dei principali dizionari esistenti. Si applicherà poi a questo lemmario un algoritmo di flessione capace di generare tutte le forme possibili secondo il sistema linguistico dato. Se invece si deve lemmatizzare un corpus per la cui lingua non esiste ancora un dizionario attendibile, si procede per gradi. Le forme del primo testo del corpus saranno lemmatizzate a mano. Il secondo testo sarà lemmatizzato adoperando un LA contenente solo le forme trovate nel primo: il LA lemmatizzerà così automaticamente solo le forme del secondo testo che comparivano anche nel primo, mentre quelle che sono nuove rispetto al primo saranno lemmatizzate a mano, e poi saranno inserite nel LA. Così il LA adoperato per lemmatizzare il terzo testo sarà costituito dall'unione delle forme dei due testi precedenti, e così via.

cale, la polisemia, la scomposizione in morfemi, gli eventuali registri di uso, le costruzioni sintattiche. Queste informazioni sono espresse in alcuni casi con un linguaggio completamente formalizzato (per esempio le costruzioni), mentre in altri casi la definizione è costituita da una parte formalizzata e da un'altra in linguaggio naturale, come nella esplicitazione delle polisemie. Un algoritmo di flessione ha generato automaticamente circa 1.000.000 di forme, che sono registrate sia secondo l'ortografia corrente sia secondo l'alfabeto fonetico. Si veda in proposito Zampolli (1973 b).

#### 1.2.3.2. Consultazione di un LA.

L'algoritmo di consultazione di un LA così organizzato è estremamente semplice. Come si vede nell'organigramma di Fig. 1, in ingresso vengono posti il LA (1) e il nastro delle concordanze per forma (2) <sup>(14)</sup>.

L'algoritmo di consultazione legge una sola parola alla volta, dal nastro 2, e la confronta con le forme grafiche che compongono il LA. Si possono verificare 3 diverse condizioni.

a) *La parola è identica a una forma che nel LA figura come univoca.*

Il programma la lemmatizza automaticamente, cioè la ricopia sul nastro di uscita (5), accompagnata, oltre che dal contesto e dal riferimento, dalla analisi che nel LA è associata alla forma in questione.

b) *Nel LA non compare nessuna forma identica alla parola cercata.*

Il programma la scrive con i relativi contesti e riferimenti sul nastro di uscita (3) che, verrà adoperato per stampare le concordanze di tutte le *parole nuove*, e cioè presenti nel testo ma assenti nel LA. I lemmi (e le altre eventuali informazioni) vengono scritti a mano in questo elenco e di qui vengono perforati su schede (6) che servono sia per lemmatizzare le parole nuove (nastro 10) sia per inserire le nuove voci nel LA (LA integrato con forme nuove, nastro 9).

c) *La parola è identica a una forma che nel LA figura come omografa.*

Salvo quanto dirò in seguito a proposito delle possibilità di distinguere automaticamente gli omografi, la loro lemmatizzazione deve essere fatta a mano. Tutte le parole omografe possibili <sup>(15)</sup> del testo vengono ricopiate sul

(14) Come si ricorderà, tale nastro contiene le parole del testo ordinate alfabeticamente, ciascuna registrata in un proprio record e accompagnata da riferimento e contesto.

(15) Nella terminologia degli spogli è invalsa, forse impropriamente, l'abitudine di distinguere tra omografia *possibile* nel sistema e omografia *attuale* in un corpus. Il primo termine indica forme grafiche che, secondo i criteri di lemmatizzazione prescelti, possono appartenere a più lemmi nel sistema linguistico cui il corpus viene riferito. Il secondo termine designa forme che, oltreché possibili nel sistema, sono anche di fatto presenti nel corpus come realizzazioni di almeno due lemmi diversi. Naturalmente, il fatto che una forma omografa possibile non sia omografa attuale in un corpus può essere constatato solo esaminandone tutte le occorrenze, e si deve indicare con opportune annotazioni all'utente dello spoglio quale, tra le diverse possibilità, è quella di fatto realizzata.

nastro di uscita (4), accompagnate ciascuna dal proprio contesto e riferimento, nonché dalle due o più analisi proposte dal LA. Il nastro 4 servirà per stampare le concordanze delle forme omografe. Accanto ad ogni forma il calcolatore stamperà anche le analisi possibili. Il linguista dovrà leggere i contesti, per assegnare ciascuna occorrenza della forma all'una o all'altra delle analisi possibili. L'analisi assegnata verrà trasferita, per mezzo di schede (8) su nastro (11: contiene le parole omografe del testo, lemmatizzate). Se l'esame dei contesti di una forma rivela una analisi non compresa tra quelle già previste per tale forma dal LA, tale analisi deve essere aggiunta al LA, per mezzo di schede (7).

### 1.2.3.3. *Obiezioni all'impiego di un LA per la lemmatizzazione.*

Com'è noto, i primi LA furono messi a punto per la traduzione automatica dal russo all'inglese e viceversa, a partire dal 1946; agli inizi, ci fu anzi chi ritenne che un buon LA fosse condizione necessaria e sufficiente per tradurre automaticamente.

I primi testi di linguistica computazionale dedicano grande spazio ai sistemi per compilare e per consultare un LA<sup>(16)</sup>. Tutti questi studi hanno prodotto sistemi molto noti e affinati per la gestione e la consultazione di un LA.

Non deve però stupire se, nonostante questo, pongo l'impiego dei LA tra i possibili sviluppi della lessicografia assistita dai calcolatori. In realtà sono pochissimi le imprese lessicografiche e in genere gli autori di spogli lessicali e statistici che abbiano adottato un LA per rendere automatica o almeno semiautomatica la lemmatizzazione.

Molti obiettono che la lemmatizzazione per mezzo di un LA secondo la procedura abituale espone al rischio di gravi errori. Il più grave sarebbe quello che una forma di fatto omografa nel sistema linguistico cui il testo da lemmatizzare appartiene, sia stata invece inserita come univoca nel LA, per errore o per insufficiente conoscenza della lingua. In questo caso il calcolatore lemmatizzerebbe direttamente come univoche le occorrenze di tale forma nel testo, senza sottometerle all'esame del linguista. Più lo strato di lingua studiato si estende diacronicamente, più si arretra nel tempo verso sistemi che non sono completamente presenti alla competenza di chi compila il LA, più questo rischio è grave<sup>(17)</sup>.

(16) Si vedano per esempio gli articoli di Locke e Booth (1955), e, per una bibliografia di questo periodo, Mounin (1964).

(17) Una forma può essere omessa nel LA o perché è stato omesso il lemma corrispondente, o perché, pur essendo presente il lemma, l'algoritmo di flessione non prevede regole capaci di generare tale forma. In questo caso, non del tutto improbabile quando si ha che fare con stati di lingua non del tutto consolidati o che includono un grande numero di varietà regionali e simili (si pensi per esempio alle varianti della lingua italiana delle origini), conviene, se possibile, adottare la consultazione per segmentazione, che permette di arricchire e complicare le regole del componente morfologico e morfonologico, aumentando così il numero delle forme potenzialmente generabili dal sistema.

Tuttavia l'ostacolo è superabile se, anziché considerare definitiva la lemmatizzazione delle forme univoche secondo il LA, se ne stampano i contesti lemmatizzati affinché il linguista controlli i lemmi attribuiti dal LA (cfr. n. 12 in Fig. 1).

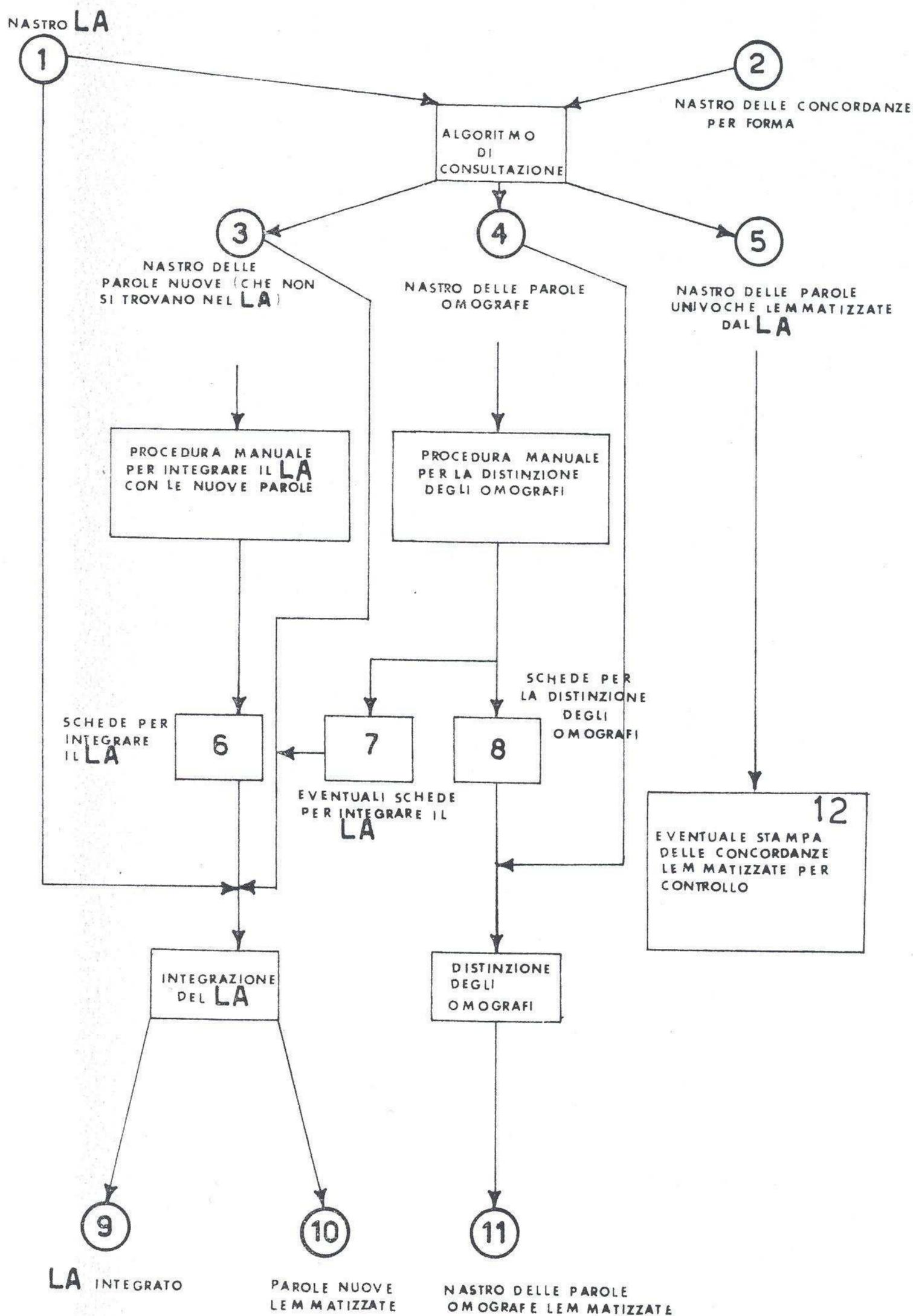


Fig. 1.

Appositi esperimenti hanno mostrato che anche in questo caso si risparmierebbe sempre un tempo notevole (il 70% circa) rispetto alla procedura di lemmatizzazione manuale.

Un'altra obiezione è quella di coloro i quali affermano che si dovrebbe creare un LA diverso per ciascuno spoglio, perché a seconda della natura del testo o degli scopi, l'autore dello spoglio può adottare criteri linguistici diversi nella lemmatizzazione. Divergenze di questo tipo certamente sono giustificabili in alcuni casi anche se, a mio avviso, è evidentemente auspicabile che i criteri di lemmatizzazione siano comuni per il maggior numero possibile di testi: si pensi per esempio alla necessità di poter comparare tra loro i dati sulle frequenze lessicali di testi spogliati da ricercatori diversi.

Tuttavia, non è vero che un LA debba necessariamente produrre un solo tipo di lemmatizzazione. Entro certi limiti, esso permette diversi tipi di lemmatizzazione, e ne assicura nel contempo la comparabilità.

Evitando il gran numero di trascrizioni necessarie nella procedura artigianale, l'uso di un LA diminuisce notevolmente il rischio di errori casuali. Nello stesso tempo funge da sistema di riferimento che assicura coerenza di comportamento nei casi nei quali la presenza di alternative sistematiche richiede che la formulazione del lemma avvenga secondo regole prefissate. Il LA funziona come registrazione delle decisioni prese e, potendo essere stampato rapidissimamente, o consultato tramite video, fornisce a ogni istante un quadro completo delle forme esaminate e dei trattamenti adottati. Accanto a questa funzione normalizzatrice, che è tanto più importante quanto più è esteso il corpus sottoposto a spoglio e quanto più numerosi e distribuiti nel tempo sono i lemmatizzatori, il LA esercita un insieme di funzioni collaterali; si pensi per es. alla possibilità, suggerita tra gli altri da C. A. Mastrelli per l'Accademia della Crusca, di retrodatare automaticamente le forme nel corpus scelto.

Un'altra obiezione riguarda il costo e il tempo richiesti dalla compilazione di un LA. Per es. per completare la parte del LAI che serve per la lemmatizzazione di testi italiani contemporanei, hanno lavorato 7 persone per 2 anni. Si deve però considerare che questo tempo equivale al tempo che il LAI permetterà di risparmiare nella lemmatizzazione di 2 milioni di parole, appena un ventesimo delle parole presenti nella biblioteca elettronica del CNUCE.

Occorre soprattutto considerare che il LA, e quindi il tempo speso per compilarlo, sono utili non solo per la lemmatizzazione, ma anche per numerose ricerche linguistiche, alcune delle quali sono premesse necessarie per lo sviluppo stesso della lessicologia e della lessicografia, sia in quanto discipline autonome, sia in quanto assistite dal calcolatore. Accanto a statistiche sulle unità dei diversi componenti che costituiscono il sistema linguistico italiano<sup>(18)</sup>,

(18) Si vedano, per esempio a proposito delle statistiche fonologiche, Troubeckoy (1939 cap. 7) Martinet (1960, n. 3.38), Hockett (1955 e 1966) Kramsky (1964), Tagliavini (1968) Muljacic (1969, p. 306) Zampolli (1968 e 1973 b), Saramandu (1966), Guiraud (1954 e 1967), Dubois (1964).

il LAI permette la costruzione di modelli dei meccanismi operanti nella organizzazione del lessico <sup>(19)</sup> e di una banca di dati linguistici <sup>(20)</sup>.

1.2.3.5. *Distinzione automatica degli omografi.*

Risulta evidente dalla procedura esposta al n. 1.2.3.2. che per completare la automazione della lemmatizzazione occorrerebbe poter distinguere automaticamente gli omografi. Questo problema è stato ed è affrontato da tempo nei progetti di traduzione automatica, mentre conosco solo pochi tentativi da parte di imprese lessicografiche. A Nancy l'èquipe del *Trésor* aveva iniziato a studiare degli algoritmi per gli omografi di altissima frequenza, e lo stesso stiamo facendo a Pisa. Analoghi tentativi sono in corso presso la 'Accademia Reale Spagnola' (Del Campo *et al.*, 1974).

Per l'omografia funzionale (del tipo sostantivo-verbo: *faccia*), si propone un *parser* sintattico.

Il principio è evidente. Il LA associa a un omografo di questo tipo, non una sola analisi grammaticale, ma più analisi distinte, una per ciascuna diversa funzione sintattica che l'omografo potrebbe esercitare. Se il *parser* ha successo e assegna un solo indicatore sintagmatico alla frase, allora il calcolatore potrà risolvere automaticamente l'omografia scegliendo la descrizione grammaticale dell'omografo che ha permesso al *parser* di avere successo.

Per l'omografia di tipo radicale, del tipo *mozzo* (della ruota e della nave), che non comporta diverse analisi grammaticali, si proponeva fino al 1966-67 un *parser* a livello semantico. Si pensava per esempio di classificare tutte le voci presenti nel LA secondo categorie o comunque componenti semantiche, e di formulare delle regole specificanti la possibilità o impossibilità di relazione, di selezione, di cooccorrenza tra le diverse categorie semantiche, almeno in certe posizioni della struttura sintattica, con il proposito, più o meno esplicito, di costruire un modello semantico globale, un reticolo di rapporti collegante, al limite, tutte le parole del lessico.

Oggi la situazione è cambiata e si tende a non separare i due momenti del *parser*. Lo stato delle procedure di analisi sintattica e semantica è descritto in Zampolli (1974 e), al quale qui rinvio. Qui basterà osservare che i sistemi più recenti citati in tale rassegna, se da un lato riescono a trattare una grande varietà di strutture sintattiche complesse, dall'altro sono limitati di solito a un sottoinsieme lessicale molto ristretto, e il loro funzionamento presuppone che gli enunciati che vengono sottoposti ad analisi si riferiscano a un argomento preciso noto a priori, a una parte di realtà ben delimitata, la cui conoscenza è data al sistema come una rappresentazione formale esplicita.

Le occasioni di incontro tra lessicografi e lessicologi, impegnati nello spoglio di vasti corpora di testi, e linguisti computazionali impegnati nello studio dei cosiddetti sistemi integrati di analisi linguistica sono mancate anche nel pas-

(19) Si veda per esempio la ricerca sui suffissi in Zampolli (1973 e).

(20) Si veda più avanti il capitolo sulla banca dei dati.

sato più recente. Negli ultimi tempi si è cercato deliberatamente di creare delle occasioni di incontro e di discussione, da parte di ricercatori che, come lo scrivente, avendo esperienza in entrambi i campi, sono convinti della possibilità e necessità di scambi metodologici.

Nel corso della Scuola Estiva 1972 e della ICCL 1973 già citate, gli autori dei sistemi di *parser* più recenti, interrogati da lessicografi e lessicologi sulla possibilità di estendere i loro modelli fino a includere sottoinsiemi sempre più vasti, e al limite tutta la lingua di uso comune, hanno espresso giudizi contrastanti sulla possibilità teorica di tale estensione, ma hanno tutti concordato sulla impossibilità pratica di realizzarla concretamente, almeno per numerosi decenni. Ciò perché, a differenza dei sistemi di *language understanding*, che si occupano di sottoinsiemi di lingua naturale estremamente ridotti, e per lo più in certa misura già formalizzati, i lessicografi sottopongono a spoglio corpora di testi distribuiti ampiamente sia diacronicamente sia sincronicamente (regione, argomento, registro d'uso, genere letterario, ecc.). Si deve quindi concludere che la soluzione automatica di tutte le omografie di un testo è almeno oggi una meta utopistica. Ma non si deve escludere né rinunciare ad automatizzare in parte la soluzione dell'omografia. Anche se, come vedremo, è raro che venga studiata o almeno resa nota l'efficienza dei sistemi di *parser*, possediamo comunque alcuni dati indicativi, secondo i quali su testi di lingua inglese, scelti senza particolari restrizioni, sarebbe possibile risolvere automaticamente, incorporando alcuni dei ritrovati più recenti dei sistemi integrati, il 60% delle omografie. Con algoritmi ad hoc, che, pur incorporando tali ritrovati, non mirino a riconoscere la struttura completa della frase, ma si proponano di risolvere soprattutto l'omografia funzionale la quale, come è noto, è generata per la maggior parte da poche parole di altissima frequenza, sembrerebbe possibile raggiungere e in talune lingue sorpassare l'80%. Risultato apprezzabilissimo perché, a differenza dei sistemi integrati ove importa il processo di analisi per se stesso, negli spogli importa essenzialmente ridurre il lavoro umano.

Probabilmente i primi passi concreti dovrebbero essere mossi verso una interazione uomo-macchina nello scioglimento della omografia, in attesa di raggiungere uno 'stato dell'arte' nel quale il programma dovrebbe eseguire il parsing delle frasi per le quali è adeguato, e dovrebbe richiedere la collaborazione del linguista per quelle che sono al di fuori delle sue capacità. Questo colloquio uomo-macchina richiede l'uso di un terminale video. Alcune dimostrazioni di questo colloquio le abbiamo organizzate con B. Quemada per i partecipanti alla Scuola Estiva nel 1972.

Sullo schermo appare la forma da analizzare. Accanto ad essa appaiono le diverse proposte di analisi fornite dal LA e dal *parser*, numerate progressivamente.

Al di sotto di queste informazioni, il ricercatore può far apparire, uno dopo l'altro, i contesti della forma. Se il *parser* è riuscito ad analizzare alcune occorrenze, il lemmatizzatore controlla l'esattezza dell'analisi scelta. Altrimenti l'opera egli stesso, per mezzo della tastiera o del *light-pen* scegliendo

tra le analisi proposte dal LA. Nel caso che nessuna di queste convenga ad una certa occorrenza, egli aggiungerà la nuova analisi a quelle già esistenti nel LA. Evidentemente quest'ultimo caso è equivalente, dal punto di vista della procedura, al caso di un testo da analizzare per la cui lingua non esiste un LA: mano a mano che l'analisi procede, si forma un LA che viene arricchito dagli spogli successivi.

#### 1.2.3.6. *Selezione dei materiali lessicali.*

Il tema della selezione domina il processo di compilazione dei dizionari e in particolare dei dizionari storici <sup>(21)</sup>.

Per incominciare, i compilatori devono scegliere un corpus di testi, che possa essere considerato rappresentativo, ai fini del dizionario, del corpus totale dei testi teoricamente disponibili nella lingua in questione <sup>(22)</sup>. Non si può negare l'opportunità di una sorta di *feedback*. Mano a mano che lo spoglio, procedendo, fornisce nuove informazioni sulla struttura del campione in esame, dovrebbe essere possibile modificare il corpus stesso, per esempio riducendo il numero dei testi relativi a un sottoinsieme di lingua o a classi di fatti linguistici sufficientemente documentati; oppure, viceversa, introducendo testi che si suppone contengano fenomeni non ancora apparsi, ecc. Questo tipo di *feedback*, fino ad oggi pochissimo realizzato, richiede una organizzazione particolare dello spoglio che consenta di gestire un archivio costantemente aggiornato e rapidamente accessibile in tutte le sue parti. Le caratteristiche di un siffatto archivio sono quelle proprie della cosiddetta *banque de mots*, di cui parleremo al numero seguente, organizzata attraverso un LA.

(21) Il dizionario storico tradizionale, del tipo rappresentato dall'*Oxford English Dictionary*, dalle diverse edizioni del *Vocabolario della Crusca* e da altri dizionari simili per molte lingue nazionali, può essere considerato essenzialmente come l'insieme ordinato di numerose citazioni selezionate da un vasto corpus di testi relativi alla lingua e al periodo abbracciati dal dizionario. Le citazioni che accompagnano ogni lemma sono riunite dal redattore per mostrare le forme, gli usi e le collocazioni di ciascuna parola e le distribuzioni di ciascuna di queste caratteristiche nel tempo, nello spazio, nei livelli di stile, ecc. « From this point of view the definitions serve merely as sign-posts or markers to the group of citations which follow them, and the primary function of the dictionary is to present the selection of citations and only secondarily to provide a list of definitions or 'meanings' ». (J. A. Aitken, 1972). Di conseguenza dizionari di questo tipo, che presentano raccolte di parecchie centinaia di migliaia di esempi, sono principalmente il frutto di un procedimento di selezioni successive.

(22) Per esempio: J. A. Aitken ha calcolato per che il DOST (Dictionary of the Older Scottish Tongue) il totale delle parole nei testi disponibili di scozzese antico sia dell'ordine di 1.000 milioni. In pratica solo un quinto di queste, circa 200 milioni, (pari a 3.000 volumi) possono essere incluse nel corpus. La scelta del corpus campione avviene di solito sulla base di giudizi sull'importanza e sulla rappresentività dei testi, o in seguito a fattori contingenti, per esempio la disponibilità di edizioni critiche. Sono molto rari i casi nei quali è possibile spogliare l'intero corpus di testi disponibili: si tratta per lo più di sottoinsiemi di lingua ben delineati o di lingue morte delle quali sono sopravvissuti relativamente pochi testi.

Il secondo processo di selezione consiste nello scegliere da corpus gli esempi che devono confluire nell'archivio.

I lessicografi sono d'accordo, in genere, con il calcolo di Aitken secondo il quale un esperto lessicografo può esaminare, redigendo un articolo di dizionario, circa 10.000 schede contesto all'anno. A questa velocità elaborare una documentazione composta da 10.000.000 di citazioni richiederebbe 100 collaboratori per 10 anni. In pratica gli archivi dei grandi dizionari storici basati su spogli manuali negli ultimi 150 anni comprendevano da un mezzo milione a 10 milioni di occorrenze, e ciononostante i preventivi di tempi sono stati quasi sempre superati di molte decine di anni. L'impiego del calcolatore, che permette di raccogliere in un tempo relativamente breve un numero di esempi molto superiore, rende, se possibile, più drammatico il problema di scegliere le occorrenze particolarmente « interessanti » da inserire nell'archivio <sup>(23)</sup>.

L'Accademia della Crusca e il CNUCE hanno messo a punto una procedura per la composizione di un archivio lessicale che è, a tutt'oggi, da un punto di vista globale, tra le più economiche e perfezionate <sup>(24)</sup>. Lo spoglio elettronico prevede la registrazione integrale del testo in *machine readable form* e la produzione delle concordanze di tutte le forme del testo. Una equipe di lemmatizzatori legge le concordanze per forma lemmatizzandole, e nel contempo contrassegnando le occorrenze ritenute degne di entrare nell'archivio. Alcuni esperimenti per automatizzare almeno in parte questa selezione sono stati compiuti al C.N.U.C.E. e presso il T.L.F. di Nancy. Essi si fondano soprattutto su metodi statistici e i risultati ottenuti sono ancora molto discutibili. Nel corso della Scuola Estiva di Pisa del 1972 si è ampiamente discusso dei mezzi atti ad abbreviare queste operazioni di scelta e a fondarle, se possibile, su criteri non esclusivamente dipendenti dal giudizio individuale del lemmatizzatore. L'idea di base prevede l'interazione del lessicografo con il programma, per mezzo di uno schermo video. Il calcolatore assisterebbe attivamente il lessicografo: nel procedimento di lemmatizzazione, mediante la consultazione di un dizionario di macchina; nel processo di selezione, mediante il riconoscimento automatico nel testo di alcune strutture descritte in precedenza come fattori rilevanti per l'accettazione o la esclusione degli esempi, e mediante il continuo controllo statistico dell'accumularsi degli elementi scelti.

Il terzo processo di selezione appartiene alla fase di redazione vera e propria del dizionario. I redattori devono analizzare i materiali raccolti nello archivio, scegliere gli esempi più rappresentativi ed organizzarli opportuna-

(23) La forma tradizionale dell'archivio è quella di un insieme di schede-contesto; ad ogni parola scelta nei testi spogliati viene intestata una scheda, nella quale si ricopia anche il riferimento al luogo del testo da cui proviene e il contesto nel quale appare. Le schede vengono ordinate, di solito, secondo l'ordine alfabetico delle intestazioni (forme, lemmi o lessemi) e, all'interno di una stessa intestazione, in ordine cronologico, di autore, di opera, ecc.

(24) Essa prevede tre tipi di spoglio, distinti dalla densità: lo spoglio fitto viene realizzato con mezzi elettronici; lo spoglio medio con procedure xerografiche; lo spoglio rado con la tradizionale schedatura a mano.

mente per illustrare ciascuna accezione, la sua distribuzione ed evoluzione nel tempo, nello spazio, ecc., scrivere la definizione e alla fine seguire fino alla stampa il risultato di tutte le operazioni. Si può immaginare un procedimento di redazione assistita dal calcolatore. Il LA fornirebbe innanzitutto un primo schema della struttura dell'articolo, eventualmente modificabile in base alle nuove evidenze emergenti dall'esame dell'archivio.

Il redattore potrebbe far comparire in sequenza sul video gli esempi, e, per mezzo di semplici comandi sulla tastiera o con il *light-pen*, eseguire in una stessa fase le seguenti operazioni:

indicare con dei numeri il raggruppamento degli esempi in classi;

chiedere al calcolatore di ripresentare gli esempi raggruppati secondo le classi prescelte;

indicare gli esempi che devono confluire nel dizionario, eventualmente tagliando nel modo più opportuno quelli troppo lunghi;

inserire in testa ad ogni gruppo la definizione o, più generalmente, i classificatori pertinenti;

inviare i materiali così elaborati (definizioni e esempi) ad un processo automatico di stampa, per es. la fotocomposizione.

Uno schema di questo tipo è stato discusso a Pisa (Aitken, 1972; Bailey, 1972; Quemada, 1972; Zampolli, 1974 b) e la Sezione Linguistica del CNUCE sta sviluppando una procedura interattiva che completa un primo esperimento presentato nel 1972 alla Scuola Estiva <sup>(25)</sup>.

Si può giustamente obiettare a questo schema che non sarebbe possibile operare, sullo schermo, l'esame sinottico dei diversi esempi, cosa che il lessicografo è abituato a fare « sparpagliando » le schede-contesto sulla scrivania. Alcuni hanno proposto di rimediare mettendo due video in parallelo, su uno dei quali sarebbe fisso il contesto in esame, mentre sull'altro verrebbero fatti scorrere dal ricercatore gli altri contesti della stessa forma. Altri hanno proposto di ricorrere a tecniche particolari, che combinano assieme un terminale video e le microfiches. Sembra certo che occorrerà ancora un serio e forse lungo lavoro di sperimentazione in stretta collaborazione tra linguisti computazionali e lessicografi. In ogni caso questa è certamente una delle più interessanti prospettive di sviluppo, che si inserisce naturalmente nel progetto più generale di una *banque de mots*.

(25) Il ricercatore può richiedere al calcolatore di raggruppare immediatamente i contesti che hanno ricevuto la stessa analisi, e di far scorrere ciascun gruppo sullo schermo, preceduto dalla analisi e da altri eventuali commenti da lui inseriti in precedenza. Nello stesso tempo è possibile operare, per es. con l'opportuno uso del *light-pen*, una scelta dei contesti da inviare all'articolo del dizionario, ed eventualmente « tagliare » i contesti troppo lunghi nel modo più opportuno. Basterà infatti per esempio appoggiare il *light-pen* sopra il primo carattere di una parola e sotto l'ultimo carattere di un'altra per comandare al calcolatore di eliminare dal contesto tutte le parole comprese tra i due punti segnati, sostituendole per esempio con dei puntini tra parentesi.

### 1.3. LA BANQUE DE MOTS.

Alla Scuola di Pisa si sono delineate due diverse interpretazioni del termine *banque de mots* o archivio lessicale.

Da un lato alcuni, come Bahr (1972) l'hanno inteso come una specie di LA, nel quale ogni parola è accompagnata da tutte le informazioni fonologiche, morfologiche, sintattiche, semantiche note. Quemada ed io l'intendevamo come un vero e proprio archivio di parole estratte da corpus di testi analizzati, elaborati, e non stampati, ma conservati per il calcolatore e messi a disposizione dei linguisti e dei lessicografi mediante le moderne tecniche conversazionali di interazione uomo-macchina.

Penso che una *banque de mots* o archivio lessicale debba contenere entrambi questi elementi, il LA e i testi, e che, poiché l'archivio è pensato come dinamico, sia il LA sia il corpus debbano avere caratteristiche dinamiche.

Il LA può costituire la registrazione ove si depositano le conoscenze acquisite da linguisti e lessicografi sul lessico e sulle sue relazioni funzionali con gli altri componenti del sistema linguistico. Il LA fornisce così un inventario delle unità del sistema lessicale e delle loro proprietà, le quali fungono da classificatori nei confronti delle unità del testo. L'unità, linguisticamente definita e classificata del LA, è la chiave attraverso la quale normalmente si accede alle unità corrispondenti del corpus.

L'indexamento del corpus richiede il riconoscimento e la classificazione di tali unità nel quadro di una teoria linguistica. Le teorie linguistiche, soprattutto oggi, sembrano mutare piuttosto rapidamente. È inevitabile che chi spoglia corpus di notevoli dimensioni si ponga il problema se, in queste condizioni, i risultati del proprio lavoro siano così legati a una teoria linguistica specifica, da risultare scarsamente utilizzabili o addirittura non interpretabili in una teoria diversa.

La risposta a questo interrogativo, che si pone in maniera tanto più drammatica quanto più ampio è il corpus e quindi proporzionalmente più costoso lo spoglio, non è facile, come si è potuto constatare nel corso delle vivaci e talora aspre discussioni su questo argomento che hanno avuto luogo al *Colloque sur l'indexation maximale* di Strasburgo (1972).

In ogni caso essa non può consistere in un *si* o in un *no* netti, ma deve invece prendere la forma di procedure cautelative che assicurino, nei limiti del possibile, la dinamicità dell'analisi linguistica del corpus.

Tra l'atteggiamento estremo di chi sostiene che l'analisi di un corpus dipende completamente da una teoria e non ha senso se non è spinta fino al livello di dettaglio massimo previsto dalla teoria, e l'altro estremo di chi, dalle stesse considerazioni, conclude che l'indexamento deve arrestarsi alle unità definite al livello grafemico, lasciando all'utente degli indici ogni analisi, si pongono a nostro avviso dei procedimenti intermedi. Questi processi sembrano poggiare su una convinzione, talora non esplicitamente asserita, che oltre alle unità puramente grafemiche sia possibile individuare delle unità























































































































