

# La digitalizzazione del *GDLI*: un approccio linguistico per la corretta acquisizione del testo?

Eva Sassolini, Marco Biffi, Francesca De Blasi,  
Elisa Guadagnini, Simonetta Montemagni

<https://aiucd2021.labcd.unipi.it/wp-content/uploads/2021/01/a067.pdf>

# Il contesto della ricerca: il progetto TrAVaSI (Trattamento Automatico di Varietà Storiche di Italiano)

## TrAVaSI

- nasce dalla collaborazione tra l'Istituto di Linguistica Computazionale "A. Zampolli" del CNR e l'Accademia della Crusca
- è co-finanziato dalla Regione Toscana con le risorse del POR FSE 2014-2020 – Asse A Occupazione, all'interno di "GiovaniSì", il progetto regionale Toscano per l'autonomia dei giovani

## Obiettivo

- costruire risorse e strumenti per il trattamento automatico di varietà storiche della lingua italiana, al momento pressoché inesistenti nel panorama nazionale e internazionale. L'obiettivo oggetto del presente contributo consiste nell'estrazione e strutturazione della conoscenza contenuta in dizionari storici digitali, al fine di creare i presupposti per funzionalità avanzate di navigazione e interrogazione

## Il punto di partenza

- uno degli strumenti realizzati dall'Accademia della Crusca all'interno di altri progetti, ovvero la versione elettronica del *Grande Dizionario della lingua italiana* (GDLI), un prototipo della quale è attualmente consultabile negli *Scaffali digitali* del Sito Web dell'Accademia

**La strutturazione e marcatura del GDLI nella sua versione informatizzata è oggetto del presente contributo**

# GDLI e [www.gdli.it](http://www.gdli.it)

- Il *GDLI* è un dizionario storico della lingua italiana, l'unico a coprire, per quanto con alcuni limiti, l'intero arco cronologico dalle origini fino a oggi. È pertanto uno strumento fondamentale per lo studio dell'italiano in diacronia.
- L'opera è suddivisa in 21 volumi, usciti tra il 1961 e il 2002, più 2 supplementi (usciti nel 2004 e nel 2008) e un volume di indici delle abbreviature delle fonti citate.
- L'attuale versione elettronica ([www.gdli.it](http://www.gdli.it)), resa disponibile in rete dall'Accademia della Crusca in seguito a un accordo con la casa editrice UTET, consente un'interrogazione *full text* per forme del testo ottenuto con applicazione di OCR senza nessuna collazione di controllo. L'approdo finale di una ricerca è il testo elettronico agganciato alla riproduzione in facsimile delle pagine.
- Vista la grande importanza dello strumento e le enormi potenzialità di ricerca seppure in una versione con errori di riconoscimento, l'Accademia della Crusca ha deciso di mettere a disposizione degli studiosi di lingua italiana anche questa versione provvisoria e ha avviato, di concerto con l'ILC, il progetto di implementazione che qui presentiamo, volto anche all'individuazione di procedure, metodi e strumenti più generalmente applicabili alla lessicografia elettronica, diretta (per la costruzione di dizionari) e indiretta (per la trasformazione di dizionari cartacei preesistenti).

# GDLI: il dizionario originale

- **Caratteristiche fisiche:** il dizionario è stampato su carta non bianca e semilucida e conta 23 volumi di grande formato, per un totale di oltre 25.000 pagine
- **Caratteristiche editoriali:** ogni pagina del dizionario è divisa in tre colonne, il font è fitto e schiacciato, la dimensione carattere può arrivare fino a 6.5 punti, l'interlinea è molto piccola
- **Caratteristiche lessicografiche:** il testo è fortemente strutturato (ogni voce può avere da tre a otto campi), sono presenti diverse varietà di italiano, si riscontrano molte abbreviazioni e caratteri speciali



Il punto di partenza

Il processo di strutturazione dei dati

L'approccio linguistico

Il supporto alla correzione manuale

# Il dizionario originale

DIVULGAZIONE

2. Traduttore.  
*Monte*, 3-112: Dopo aver consultato un suo calepino d'indirizzi s'era fatto indicare da un droghiere la casa da lui desiderata... Abbiava dunque il Ponzio Macchi, il più instancabile e forse il più sottile dei suoi divulgatori stranieri?

3. Ant. Che celebra, che esalta.  
*Cecchi*, 5-24 (1-89): Voi siete divota e serva d'essa madama la reina e continua e chiara divulgatrice de le sue rare doti.

= Voce dotta, lat. tardo *divulgator -oris*, da *divulgare* 'divulgare'.

**Divulgazione** (ant. *divulgatione*), sf. Il divulgare, il rendere noto pubblicamente; diffusione di notizie, idee, dottrine, usanze, scritti.  
*Tasso*, II-135: Io aspetto la ricompensa di quel dispiacere che mi ha portato la divulgazione de l'opere mie così maltrattate. *Davila*, 29: Io non posso indurmi ad affermarlo sopra la sola e molte volte fallace divulgazione della fama. *Tommaso*, 1-138: Il Cioni mi manda la risposta che il Beccbi fece alla sua lettera intorno all'edizione de' novelli e de' comici. Ma io non vo' cooperare alla divulgazione di sozzure. *E. Cecchi*, 5-452: Mi ricordo certe figure che erano in voga, anni addietro, durante la prima e più entusiastica divulgazione del progresso scientifico ed industriale. *Brasconi*, II-201: Ormai non aveva più fiducia che una sola parola, pronunciata a bassa voce in quel salotto, potesse sfuggire alla divulgazione.

- Ant. In senso concreto: notizia, voce, diceria.

*Davila*, 241: Ma molto di più di questa divulgazione, universalmente creduta falsa, empirono il Re di sospetto le lettere di Monsignore di San Goart.

2. In partic., con riferimento a certa produzione artistica o letteraria: esposizione chiara e formulata in un linguaggio largamente comprensibile, di nozioni scientifiche e tecniche, di dottrine, di ideologie che divengono in tal modo accessibili a un pubblico assai più vasto della cerchia intellettuale che le ha elaborate (e può essere implicita, nel termine, una connotazione spregiata: relazione alla mancanza di originalità di tali opere e, soprattutto, a certa superficialità e approssimazione che a volte vi si riscontra).

*E. Cecchi*, 8-117: La letteratura di consumo corrente è francese: romanzi, libri gialli, periodici letterari e di varietà; collane filosofiche e psicologiche, con abbondanza di divulgazioni freudiane. *Serra*, 1-271: Questo è infine il nostro carattere più vero, nella cultura, come nell'arte: la banalità, l'imprecisione; la grossolanità delle disposizioni generiche senza la vita e il rilievo dei particolari. È una specie di divulgazione, di livellamento democratico. *Gramsci*, 6-106: Perché non è nata una letteratura di divulgazione scientifica, come in Francia e negli altri paesi?

= Voce dotta, lat. tardo *divulgatio -ōnis*, da *divulgare* 'divulgare'.

**Divulghévolméte** (*divulghévolménte*), avv. Ant. Pubblicamente.  
*Fabrizio Massimo* *vigilar*, 1-89: Et allora adì che colui era la credde morte di Ciglia e d'Italia, e che quando il fosse scoltato di quelli legani sarebbe strigimento di molte città; il quale sogno è di seguente rigone divulghévolméte.

Comp. di un *divulghévolméte*, non attento.  
**Divulsiōne**, sf. Letter. Il divulgare, lo strapare.

*Faldella*, 3-470: Il giorno dopo la catastrofe, il paesaggio si fece veder limpido, come nulla di strano, come fosse accaduto. Le spighe dei campi, che pur sentivano sul collo le carezze del vento, e le foglie dei gelhi timorosi della prossima divulsione non fecero scappare argomento del loro moritorio la tragedia del 2 giugno.

2. Medico. Dilatazione violenta, operata a mano o mediante appositi strumenti, di orifici muniti di anelli sinterici, che non provoca la lacerazione rendendone così la contrazione impossibile (ed è praticata nella cura del cardiopneumismo e delle ragadi anali).

*Famini*, IV-202: 'Divulsione'. In medicina significa dilatazione forzata... 'Divulsione del pube, del edo'.

= Voce dotta, lat. *divulsio -ōnis* (da *divellere*: cfr. D. *VELLERE*).

**Divulso**, v. **DIVELTO**.

**Divulsōre**, sm. Medico. Strumento chirurgico usato per praticare la divulsione.

= Voce dotta, deriv. da *divulso*, part. pass. di *divellere*: cfr. *DIVELLERE*.

**Dicitte**, sm. Ant. Solo nella locuz. *Avere più anni del dicitte*: essere antichissimo.

*Nomi*, 3-141: Provverbi se a resistar sarà buona / la loro murgaglia e quelle palafite / che oggi mai han più anni del dicitte.

= Dal lat. *dicitū*, 3ª pers. perf. ind. di *dicere* 'dire' (premissa: toscanamente *divitū*), frequente formula d'izio del Salmo.

**Diaccarere** (dial. *diaccararo*), tr. (*diaccare*, *chero*). Ant. Pulire un abito inacciarato.  
*S. Bernardino da Siena*, 858: Di verno infangasi i guastai il vestimento da piei, ché s'involte nel fango come fa una porca, e poi vi perde uno di a diaccararo.

= Comp. da *di-* con valore di separazione e *accare* (v.).

**Diceccolare** e deriv., v. **DIZCECCOLARE** e deriv.  
**Diziferare** (*diziziferare*, *dizizifare*), tr. (*diziferare*). Ant. e dial. Decifrare.

*G. G. Coste*, 237: Oh, questa è una cosa che la sanno tutti i nostri peccatori lassù... È possibile questo? Io lo tengo per una cosa molto inutile... Oh, lo ve la voglio diziferare or ora. *Citari*, II-196: M'ero quasi adattata nell'animo mio alle dure mie circostanze, quando a diadare questo altro esigma, m'arrivò da Napoli una lettera di D. Virginia.

= Variante dial. di *di-* (da area settentr.) di *decifrare* (v.); cfr. *di-* (v. 109).

**Diziferato** (part. pass. di *diziziferare*), agg. (*diziftrato*). Ant. e dial. Deciftrato.

*Citari*, II-93: Ecco diziferato l'enigma che il conte Frisano nella sua lettera indicato m'avea.

**Dizigodole**, agg. (plur. m. -'i). Biol. Derivante dalla fonocombinazione di una cellula-uovo diversa da quella da cui ha origine il gemello; dicoriale.

= Voce dotta, comp. da *di-* (dal gr. *di-* 'doppio') e *zigote* (v.).

**Dizionario**, agg. (plur. m. -'i) Raro. Poeta, dante, libresco.

*Bacchelli*, 2-XXII-253: In Italia, non sono mai occorse quelle revisioni e importazioni di gusto e fondazioni di nuovi canoni, che l'Indole accademica o dioniziana della letteratura francese e l'Indole troppo fresca e mallo-bilissima della tedesca hanno periodicamente prodotto.

Deriv. da *dizionario*, sul modello di *dicione*.

**Dizionario**, sm. Opera che raccoglie, accompagnata da una definizione, le parole di una lingua, nella loro totalità, sincrona o diacronica o limitatamente a particolari età della tradizione linguistica, o a particolari autori, o a particolari categorie di parole (zoologiche, barbarismi, sinonimi, tecnicismi delle varie scienze, arti, mestieri): l'ordine in cui le parole sono disposte è generalmente alfabetico, più di rado etimologico; non infrequente è l'ordine metodico (per cui le parole sono raggruppate secondo la loro affinità concettuale). - **Dizionario bilingue**: quello in cui i vocaboli e le locuzioni di una lingua sono definiti con quelli corrispondenti di un'altra lingua (ed è strumento essenziale per le traduzioni). - **Dizionario storico, geografico, biografico, degli autori e delle opere**: raccolta ordinata di nomi propri di personaggi e luoghi e di titoli di opere, accompagnati dalle notizie relative. - **Dizionario enciclopedico**: riunisce i caratteri sia del dizionario linguistico (in quanto tiene conto anche del patrimonio lessicale di una lingua), sia del dizionario storico, geografico, scientifico, tecnico.

*P. F. Giambullari*, 5-234: E egli però vero che tutte queste voci siano aranne? Certissimo, gli risposi io. E che se ne mostra? mi soggiunse egli. Ed io: I dizionari stessi caldi ed ebri che si trovano oggi stampati.

*Fiamma*, 296: Ha scritto ancora gli Avvertimenti morali sopra tutta la Bibbia; dato buon principio al Dizionario teologico. *Sordani*, 10-125: Egli primo di tutti aveva composto la grammatica della lingua giapponese e i più dotti dizionari. *Magalotti*, 1-28: In oggi... ogni mestiero ha la sua Piantina e il suo gran Dizionario... *Milizia*, II-123: Fece un corso completo d'architettura. Fu il primo a fare un dizionario di quest'arte. *Da Ponte*, XXIII-18: Era quel libro un dizionario, tedesco e italiano: a' loci indicati lessi queste tre parole: «Ich liebe Sie»; e trovai che significavano «Io amo voi». *Tramador* (v. 1): Il titolo di 'Vocabolario' o 'vocabolista' non si applica propriamente che a puri 'dizionari' di parole. 'Dizionario' è d'un significato più esteso, poiché comprende in generale non solo i 'dizionari' delle lingue, ma anche i 'dizionari' storici, e quelli delle scienze e delle arti. *Tommaso*, 3-1-141: Nel dizionario cito desidererei gioverebbe, 3-1-141: I modi schietti del popolo da quelli degli scrittori, che sono sovente come il lingua morta. *De Sanctis*, 1-27: Eccoti un vasetto. In un dizionario per categorie, alla voce vasello troverai infinite voci toscane per nominare le sue parti ed i suoi attrezzi. *Fosco*, e ne aveva appoggiato il dorso al letto inferiore della libreria.

- **Figur.** Schierz.

*Monti*, 3-6-103: Trimungio... sapeva intero a mente delle buone trame il dizionario. *Garavanti*, 3-9-162: Andando innanzi di questo passo tu devi essere arrivato presto a comporre il dizionario della lingua degli occhi.

*Vico*, 251: Ideo uniformi nate appi interi popoli tra coloro non conosciute debbon avere un motivo comune di vero. Qual'esse il dizionario mentale, da dicitte 'dicitte' a tutte le lingue articolate diverse, col quale sta concepita la storia Ideal eterna che ne dà la storia in tempo di tutte le nazioni. *G. Coste*, 1-131: Le prime parole che noi tutti impariamo... ci vengono dagli occhi nella memoria nel mezzo delle balie, poco e poco.

2. **Dizionario**, sm. Lessico; linguaggio.  
*Vico*, 251: Ideo uniformi nate appi interi popoli tra coloro non conosciute debbon avere un motivo comune di vero. Qual'esse il dizionario mentale, da dicitte 'dicitte' a tutte le lingue articolate diverse, col quale sta concepita la storia Ideal eterna che ne dà la storia in tempo di tutte le nazioni. *G. Coste*, 1-131: Le prime parole che noi tutti impariamo... ci vengono dagli occhi nella memoria nel mezzo delle balie, poco e poco.

3. **Dizionario**, sm. Ant. Pulire un abito inacciarato.

tamiami da tante bocche e penne servili, bisogna alla libertà, alla proprietà, ai diritti dell'uomo, alle leggi, ad ogni cosa insomma dar nuovi nomi. *De Sanctis*, 11-154: L'amore ha anche nella poesia la sua storia: prima passionato e rozzo, a poco a poco si fa gentile... e poi sfumando nel manierato, nel galante, nel convenzionale, e la forza del primo sfuggimento diviene un dizionario di moda. *Dossi*, 521: Né il mio dizionario or rifiuto, come altre volte, tutte quelle vivaci espressioni che danno il colore ed i muscoli ad un discorso.

3. **Dimin.** **Dizionarietto**.  
*Carducci*, III-15-116: Appunto per le belle fu dato quel dizionarietto mitologico, che dalla seconda edizione in poi accompagnò le canzonette.

- **Sprezz.** **Dizionariaccio**.  
*Mazzini*, 1-294: Le vostre dichiarazioni sul principio e le lunghe mie due lettere d'istruzioni, tonatemi insieme ad un dizionario del quale non vi gioverete mai... promettere, stipulare altro.

= Voce dotta, lat. mediev. *dicionarius* (da *dicitū -ōnis* 'modo di dire' + *nom.* di *ecclōbarius* 'vocalista' o *vocabolista*). **Vocabolario** è sinonimo di **dizionario linguistico** (ma non è sinonimo di **dizionario enciclopedico**).

**Dizionariata**, sm. e f. (plur. m. -'i). Compilatore di un dizionario.

*Roberti*, VI-2-151: Era d'uopo in primo luogo provare, ciò che non fa il dizionariata, che il lusso veramente arricchisca colui che fa lusso.

**Dizione**, sf. Il dire, discorso (anche riferito a singoli e concrete espressioni linguistiche, orali o scritte).

*Speroni*, 1-4-211: Questa terza condizione del nostro versare, che è la concretezza, è la coerenza, è la molteplicità ed identità delle lettere che entrano in essa sillaba sopra la quale casca la rima, par che sia concazione propria della lingua italiana e di questa, cioè quanto termina e significa. *Monti*, IV-262: Io ho stimato bene di non tenermi che alle varianti che inducono variazioni nella sentenza, e che migliorano la dizione. *Carducci*, III-18-15: L'infrazione, a cui la patria di Dante lasciava andare con gli ultimi Medici e si distese co' Lorenzi, corrose l'incarnato e la forza della dizione. *Baldini*, 1-200: Per quanto impegnato in una dizione largamente peneumatica, non gli sentii cadere in tutto il discorso sillaba in fallo, non ci fu lo strappo d'una consonante.

- **Gramm.** Ant. **Parti della dizione**: parti del discorso.

*Varoli*, 8-2-111: Aristotele, dove favella verso il fine delle parti della dizione e divide il nome in più specie, non intende di quel proprio che al presente favelliamo noi.

2. L'aspetto formale di un'opera letteraria, quello cioè che concerne la scelta e disposizione delle parole. - Per estensione stilistica.

*Speroni*, 1-4-158: Il divino Alighieri... [chiamò] stil tragico ogni poetica dizione, la qual senza imitare senta e tegna del grande. *Abbate*, 168: Per sé può anche intarsi in qualche tratto di magnifica dizione e di giudiziooso insinuamento; ma non deve nella secca maniera del suo fraseggiare e nella erudizione astrusa costruire esempio. *Maratori*, 5-1-121: Nella sentenza poi, o vogliamo dire ne' sentimenti e nella dizione, o sia nelle frasi e parole con cui descrivono i poeti le cose, infinitamente ancora si perfeziona la natura. *Foscolo*, XI-2-316: Condense e calore di stile. *Leopardi*, II-32: La lingua francese è poverissima. Quindi la necessità di metafore, di metaforismi, di catafore, delle mie figure di dizione. *De Sanctis*, 7-162: Questo ditto artistico, questo analogico è senso artistico, il riflette in ciò che si è un dir forma, nella dizione, nella lingua. *Gramsci*, 6-364: L'azione drammatica è fatta in questa condotta con candore, con dizione semplice e quasi scialba.

Raro. **Silabario**.

*Lepardi*, 1-1302: Così pure accade nel latino, che i più antichi [scrittori] sono i più facili, e di dizione più semplicissima e grossa.

3. **Lettera** ad alta voce o recitazione di testi letterari. - Anche (con speciale riferimento ad attori, conferenzieri, annunciatori radiofonici e televisivi): arte di leggere o recitare un testo in modo da farne percepire interamente a chi ascolta la bellezza poetica o l'efficacia persuasiva (con eccellenza di pronuncia, giusta intonazione di voce, espressività).

*Orvino*, X-19-85: Anche'io ho sognato ieri sera nel piccolo teatro dioniziano. *Garavanti*, 3-9-162: Ma che cosa? l'ultima volta... La sua dizione, il suo accento, la sua voce ripetono l'antico incanto. *Ojetti*, II-435: Vorrebbe che questi giovani, di popolo e di studio, conoscessero anche storia, letteratura, dizione, trucco, per anni, prima di cimentarsi alla ribalta. *De Sanctis*, 1-171: La sua dizione era la più brillante e tuttavia la più modesta, retta da un senso delle pause, dei silenzi, da un rispetto dell'arabesco sonoro, che si intendeva meglio sapendo che costata arione proveniva dalla danza.

In partic.: modo di articolare i suoni, pronuncia.

*C. E. Gadda*, 7-129: Noi udiamo invece il toscano serbare con amico animo e cogliere agevolmente le sfumature della musica di vocaboli più colti e più gentili. *Alfieri*, all'aire quando opportunità ne richieda: opportunità... cioè un minimo senso musicale, nella dizione spontanea.

**Divulgazione** (ant. *divulgatione*), sf. Il divulgare, il rendere noto pubblicamente; diffusione di notizie, idee, dottrine, usanze, scritti.

*Tasso*, II-135: Io aspetto la ricompensa di quel dispiacere che mi ha portato la divulgazione de l'opere mie così maltrattate. *Davila*, 29: Io non posso indurmi ad affermarlo sopra la sola e molte volte fallace divulgazione della fama. *Tommaso*, 1-138: Il Cioni mi manda la risposta che il Beccbi fece alla sua lettera intorno all'edizione de' novelli e de' comici. Ma io non vo' cooperare alla divulgazione di sozzure. *E. Cecchi*, 5-452: Mi ricordo certe figure che erano in voga, anni addietro, durante la prima e più entusiastica divulgazione del progresso scientifico ed industriale. *Brasconi*, II-201: Ormai non aveva più fiducia che una sola parola, pronunciata a bassa voce in quel salotto, potesse sfuggire alla divulgazione.

- Ant. In senso concreto: notizia, voce, diceria.

*Davila*, 241: Ma molto di più di questa divulgazione, universalmente creduta falsa, empirono il Re di sospetto le lettere di Monsignore di San Goart.

2. In partic., con riferimento a certa produzione artistica o letteraria: esposizione facile, chiara e formulata in un linguaggio largamente comprensibile, di nozioni scientifiche e tecniche, di dottrine, di ideologie che divengono in tal modo accessibili a un pubblico assai più vasto della cerchia intellettuale che le ha elaborate (e può essere implicita, nel termine, una connotazione spregiata, in relazione alla mancanza di originalità di tali opere e, soprattutto, a certa superficialità e approssimazione che a volte vi si riscontra).

*E. Cecchi*, 8-117: La letteratura di consumo corrente è francese: romanzi, libri gialli, periodici letterari e di varietà; collane filosofiche e psicologiche, con abbondanza di divulgazioni freudiane. *Serra*, 1-271: Questo è infine il nostro carattere più vero, nella cultura, come nell'arte: la banalità, l'imprecisione; la grossolanità delle disposizioni generiche senza la vita e il rilievo dei particolari. È una specie di divulgazione, di livellamento democratico. *Gramsci*, 6-106: Perché non è nata una letteratura di divulgazione scientifica, come in Francia e negli altri paesi?

= Voce dotta, lat. tardo *divulgatio -ōnis*, da *divulgare* 'divulgare'.

Il punto di partenza

Il processo di strutturazione dei dati

L'approccio linguistico

Il supporto alla correzione manuale

Una pagina del GDLI

Una voce del GDLI

# L'output del sistema di OCR: il punto di partenza

- Il testo restituito dal sistema di OCR è in formato non standard (word) e di difficile gestione
- Principali problematiche incontrate:
  - sul piano della formattazione: es. mancata individuazione dell'entrata; mancato riconoscimento dei paragrafi di definizione e/o di esempi; mancato riconoscimento dei paragrafi di commento etimologico; cattiva segmentazione delle voci, ecc.
  - sul piano testuale:

1) cattiva segmentazione	1.1) introduzione di spazi bianchi dove non ci sono
	1.2) univerbazione di parole separate
2) grafie scorrette	2.1) omissione di singoli caratteri o sequenze
	2.2) cattivo riconoscimento di singoli caratteri o di sequenze di caratteri
	2.3) cattivo riconoscimento di caratteri che produce l'inserimento di caratteri aggiuntivi

Tabella 1: tipologie di errori

Il punto di partenza

Il processo di strutturazione dei dati

L'approccio linguistico

Il supporto alla correzione manuale

# L'output del sistema di OCR: il punto di partenza

## Alcuni esempi

N.	Originale cartaceo	Testo OCR
1	<b>Amminoazobenzene</b> ( <i>aminoazobenzène</i> ), sm. Chim. Composto organico classificato tra i coloranti azoici conosciuto anche col nome di <i>giallo d'anilina</i> : cristalli gialli che si sciolgono in alcole ed etere, assai meno in acqua (usato nella colorazione di prodotti alimentari e per preparare altri coloranti).	<b>Am mi no a z ob e nz è ne</b> ( <i>aminoazobenzène</i> ), sm. Chim. Composto organico classificato tra i coloranti azoici conosciuto anche col nome di <i>giallo d'anilina</i> : cristalli gialli che si sciolgono in alcole ed etere, assai meno in acqua (usato nella colorazione di prodotti alimentari e per preparare altri coloranti).
2	<b>Assolare</b> <sup>1</sup> , tr. ( <i>assòlo</i> ). Disus. Rendere solo. - <i>Assolare una carta</i> : tenere scompagnata, nel gioco, una carta di un dato segno. = Deriv. da <i>solo</i> (v.). <b>Assolare</b> <sup>2</sup> , tr. ( <i>assòlo</i> ). Esporre al sole; rendere soleggiato. = Deriv. da <i>sole</i> (v.). <b>Assolare</b> <sup>3</sup> ( <i>assuolare</i> ), tr. ( <i>assòlo</i> o <i>assuòlo</i> ). Disporre a strati. = Deriv. da <i>suolo</i> (v.).	<b>Assolare</b> <sup>1</sup> , tr. ( <i>assòlo</i> ). Disus. Rendere solo. - <i>Assolare una carta</i> : tenere scompagnata, nel gioco, una carta di un dato segno. = Deriv. da <i>solo</i> (v.). <b>Assolare</b> <sup>2</sup> , tr. ( <i>assòlo</i> ). Esporre al sole; rendere soleggiato. = Deriv. da <i>sole</i> (v.). <b>Assolare</b> <sup>3</sup> ( <i>assuolare</i> ), tr. ( <i>assòlo</i> o <i>assuòlo</i> ). Disporre a strati. = Deriv. da <i>suplo</i> (v.).
3	<b>Ammacchiare</b> <sup>1</sup> , rifl. ( <i>m'ammacchio, t'ammacchi</i> ). Raro. Nascondersi nella macchia. <i>B. Davanzati</i> , I-136: Floro s'ammacchiò: vedendo poi presi i passi dell'uscita, s'uccise.	<b>Ammacchiare</b> <sup>1</sup> rifl. ( <i>m'ammacchio, Vammacchi</i> ). Raro. Nascondersi nella macchia. <i>B. Davanzati</i> , I-136: Floro s'ammacchiò: vedendo poi presi i passi dell'uscita, s'uccise.
4	<b>Attendista</b> , agg. e sm. e f. (plur. m. -i). Neol. Chi evita di prendere posizione (e resta in attesa degli avvenimenti, riservandosi di decidere secondo il loro svolgersi). = Fr. <i>attentiste</i> (1941), da <i>attendre</i> 'attendere'. <b>Attenditore</b> , agg. e sm. (femm. -trice). Ant. Che attende, aspetta.	= Deriv. da <i>attendere</i> . <b>Attendista</b> , agg. e sm. e f. (plur. m. -i). Neol. Chi evita di prendere posizione (e resta in attesa degli avvenimenti, riservandosi di decidere secondo il loro svolgersi). = Fr. <i>attentiste</i> (1941), da <i>attendre</i> 'attendere'. <b>Attenditore</b> , agg. e sm. (femm. -trice). Ant. Che attende, aspetta.

Problematiche sul piano del testo

Problematiche sul piano della formattazione

Il punto di partenza

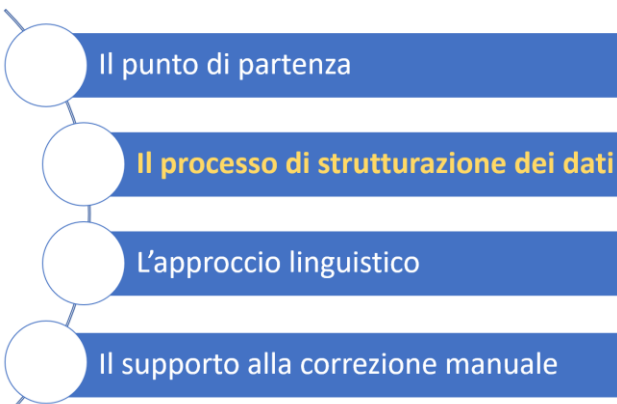
Il processo di strutturazione dei dati

L'approccio linguistico

Il supporto alla correzione manuale

# Metodi e tecnologie utilizzate per il processo di strutturazione: i dati di input

- L'output del sistema di OCR
  - la fase di acquisizione non è parte del nostro progetto
    - Impossibilità di impedire il ripetersi delle casistiche di errori più ricorrenti, individuando possibili strategie di miglioramento, cosa fattibile potendo avere un feedback dalla procedura di OCR
- La presenza di un formato dei dati digitali non standard e la presenza di errori distribuiti di vario tipo ha condizionato la scelta dell'approccio all'estrazione

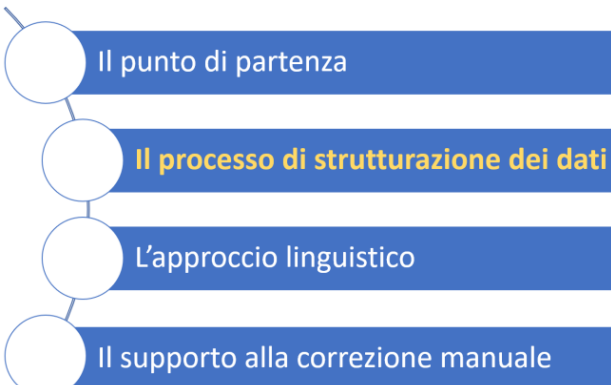




# Metodi e tecnologie utilizzate per il processo di strutturazione dei dati

## Il processo di estrazione automatica

- L'analisi dei dati ha reso necessaria una strategia di estrazione a regole personalizzate con livelli di raffinamento progressivi:
  - In una prima fase la segmentazione del testo in porzioni corrispondenti al corpo dell'intera entrata
  - successivamente si procede con il riconoscimento degli altri campi dell'entrata grazie all'identificazione delle caratteristiche distintive, tradotte in vincoli di corretta attribuzione
  - Il processo è impostato in modo incrementale



Intera voce

Abbarbagliaménto, sm. Abbagliamento intenso e improvviso della vista.

*Bencivenni [Crusca]:* Sopravviene frequente abbarbagliamento d'occhi. *Tommaseo-Rigatini, 2572:* L'abbagliamento... può essere sensazione abituale o prolungata. L'abbarbagliamento non si potrebbe immaginare continuo. *Verga, 4-275:* Un'onda di sangue al volto, un abbarbagliamento improvviso dinanzi agli occhi. *Serao, I-266:* Avevano una cert'aria inebetita, simili a chi ha assistito a un troppo lungo spettacolo musicale o coreografico, con un abbarbagliamento negli occhi e un assordamento negli orecchi. *Soffici, I-259:* L'abbarbagliamento del gran sole che al primo entrar nella stanza l'aveva lasciato come cieco. *Tozzi, III-163:* Vide, in un abbarbagliamento di sole, alcune aiuole fiorite.

2. Figur. Errore; abbaglio.

*Segneri, IV-219:* Si rivoltò [Lutero] co' Vescovi, co' Principi, co' Papi, solo perché questi lo vollero, secondo il loro debito, fare accorto de' suoi così manifesti abbarbagliamenti.

= Deriv. da *abbarbagliare*.

Singoli campi

Altro senso con citazione interna

# Metodi e tecnologie utilizzate per il processo di strutturazione dei dati: l'approccio linguistico

È possibile impiegare «intelligenza linguistica» per mettere a punto una strategia di trattamento efficace degli errori, che cerchi di ridurre al minimo la necessità di interventi manuali?

- Modularità nell'approccio:
  - Separare le strategie di individuazione degli errori da quelle di correzione automatica
  - Diversificare le strategie che interessano l'intero testo da quelle mirate al singolo campo
- Ad oggi abbiamo sviluppato procedure d'intervento sugli errori a tutto testo e nel campo lemma



Il punto di partenza

Il processo di strutturazione dei dati

**L'approccio linguistico**

Il supporto alla correzione manuale

# Metodi e tecnologie utilizzate per il processo di strutturazione dei dati: l'approccio linguistico

## Errori a tutto testo:

- Per avere un quadro generale di questi errori, abbiamo creato una «sezione *gold*»: una selezione casuale di 30 colonne, che sono state corrette manualmente
- L'analisi delle risultanze di questa selezione rappresenta lo studio preliminare su chi poggiano le strategie di correzione adottate

LIVELLO DI CORREZIONE	AREA DI TESTO INTERESSATA DALL'ERRORE	DESCRIZIONE DELL'ERRORE	STRATEGIA DI INDIVIDUAZIONE/CORREZIONE DELL'ERRORE
livello 0	singolo carattere	presenza di un carattere non alfabetico (o non appartenente all'alfabeto italiano)	• analisi degli errori della “sezione <i>gold</i> ” (sostituzioni 1:1)
livello 1	sequenza di 2 o più caratteri (ogni carattere di per sé è ammissibile)	sequenza di caratteri non ammessa in lingua italiana	• analisi degli errori della “sezione <i>gold</i> ” • “criterio fonotattico”
livello 2	parola (caratteri e sequenze di caratteri di per sé sono ammissibili)	parola errata	• analisi degli errori della “sezione <i>gold</i> ” • confronto con altri <i>corpora</i>

Il punto di partenza

Il processo di struttura

L'approccio linguistico

Il supporto alla correzione manuale

# Metodi e tecnologie utilizzate per il processo di strutturazione dei dati: l'approccio linguistico

Strategia di individuazione/correzione dell'errore: le risultanze della sezione *gold*

OCR errato	Lezione corretta
*	'
*	'
*	'
*	'
*	'
*	'
*	'

Livello 0  
Es. \* > '

OCR errato	Lezione corretta
??	zz
??	zz
??	zz
??	zz
??	zz
??	zz
??	zz
??	zz
??	zz
??	zz

Livello 1  
Es. ?? > zz

OCR errato	Lezione corretta	Passo errato	Passo corretto	Campo della voce	Fonte	vol_pa
pass,	pass.	part. pass,	part. pass.	definizione	Gold	01_20_1
pass,	pass.	part. pass,	part. pass.	definizione	Gold	01_20_1
pass,	pass.	part. pass,	part. pass.	definizione	Gold	03_488_1
pass,	pass.	part. pass,	part. pass.	definizione	Gold	06_339_3
		part. pass,	part. pass.	definizione	Gold	08_263_1
		part. pass,	part. pass.	definizione	Gold	16_1042_2

Livello 2  
Es. pass, > pass.

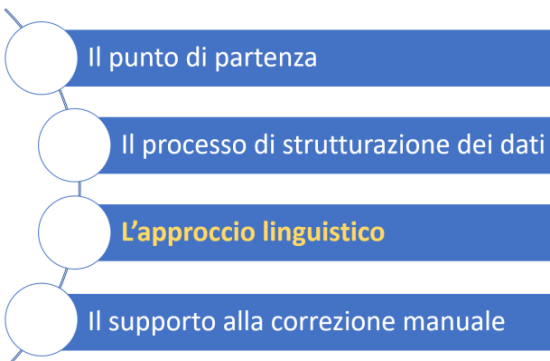
- Il punto di partenza
- Il processo di strutturazione dei dati
- L'approccio linguistico**
- Il supporto alla correzione manuale



# Metodi e tecnologie utilizzate per il processo di strutturazione dei dati: l'approccio linguistico

Strategia di individuazione/correzione dell'errore nel campo lemma

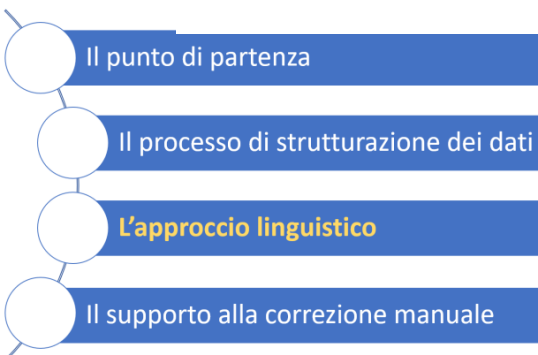
- controllo delle «testatine»:
  - confronto del lemma con l'intervallo alfabetico segnalato dalle “testatine” (vale a dire le entrate della prima voce della prima colonna e dell'ultima voce della terza colonna, segnalati nell'intestazione di ogni pagina del *GDLI*): grazie a questo confronto è possibile formulare delle ipotesi di correzione della parte iniziale del lemma – è così possibile per esempio correggere l'errato “**bfdicotomia**” in “broncotomia”.
- confronto con le entrate del *Dizionario della lingua italiana* «Tommaseo-Bellini»



# Metodi e tecnologie utilizzate per il processo di strutturazione dei dati: l'approccio linguistico

Le strategie descritte hanno mostrato la necessità di un approccio modulare

- Gli interventi correttivi si collocano in fasi e tempi diversi del processo di estrazione:
  - L'output del sistema di OCR per gli errori generalizzati (a tutto testo)
  - i testi dei campi della voce in fase di estrazione, valutando la combinazione di diverse caratteristiche
- Anche le procedure software diventano modulari e utilizzano diverse risorse linguistiche nelle fasi mirate di estrazione di singoli campi/porzioni di campo



# Realizzazione di strumenti e di risorse linguistiche a supporto della correzione manuale: valutazione dei risultati

- Una parte degli errori descritti è stata corretta automaticamente ma resta necessaria una revisione manuale post-processing
- Implementazione di strategie e metodi di supporto alla correzione e sistema di revisione e riallineamento successivo dei dati estratti
- Parallelamente alla fase di estrazione dei dati sono stati implementati meccanismi di annotazione puntuale:
  - creazione di un file di report che segnala la mancanza di campi obbligatori o se il loro ordine non è rispettato, con ubicazione precisa del punto;
  - quando è possibile impostare un'indagine più puntuale del dato, i file di report sono finalizzati al controllo delle soluzioni già inserite in fase di parsing di estrazione, così da facilitare il lavoro di revisione manuale
- Impostazione di strumenti software per supportare la correzione assistita del lemmario:
  - sulla base delle evidenze emerse nello studio della sezione *gold* e dei criteri di individuazione degli errori sopra descritti è stato stilato un preciso protocollo di correzione automatica e manuale che si avvale di un opportuno strumento software di supporto dedicato

Il punto di partenza

Il processo di strutturazione dei dati

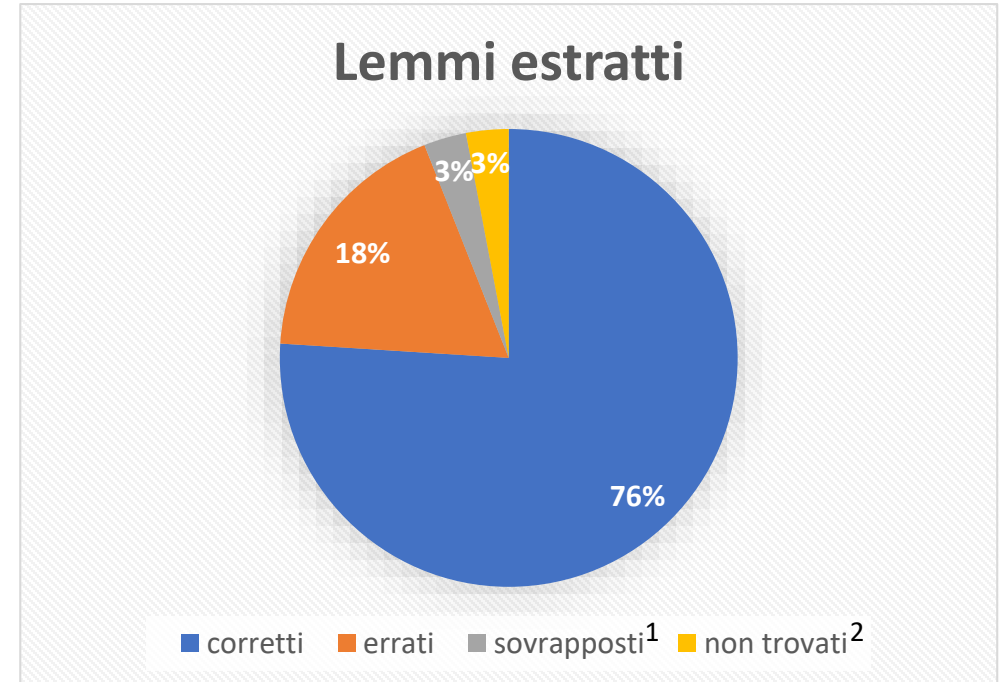
L'approccio linguistico

Il supporto alla correzione manuale



# Realizzazione di strumenti e di risorse linguistiche di supporto alla correzione manuale

- Tutti i file di report sono organizzati e gestiti da un software che assiste la correzione manuale
- Il testing dello strumento software di supporto è stato fatto impostando una prima sessione di lavoro:
  - sul primo volume
  - relativamente alla correzione del lemma
- Nell'immagine accanto è riportata la sintesi



1 **sovrapposti**: conseguenza dell'errata interpretazione dello strumento morfologico del sistema di OCR, es. «aggranchito» per «aggronchito»

2 **non trovati**: con struttura dell'entrata diversa dalla norma, es. manca la categoria grammaticale anche nell'originale

Il punto di partenza

Il processo di strutturazione dei dati

L'approccio linguistico

Il supporto alla correzione manuale

# Conclusioni

- Il lavoro già oggi rappresenta un valido prototipo di intervento in situazioni in cui la mole e la qualità dei dati da elaborare rendono necessaria la messa a punto di procedure *ad hoc*
- Il primo scopo del lavoro è portare l'output del sistema di OCR in un formato standard di rappresentazione, che consenta la strutturazione del vocabolario e quindi moltiplichi le sue possibilità di consultazione e interoperabilità
- In prospettiva è quindi pensabile utilizzare altri strumenti software di analisi per affinare ulteriormente la marcatura dei dati