

Digital tools in linguistic studies: Unlocking language heritage

Abstract¹

Over the years, digital tools have played a pivotal role in advancing the study of language. They not only enhance our understanding but also open up new research avenues. One significant aspect of this digital revolution is the digitisation of book heritage, a process that not only benefits specialists but also makes valuable materials accessible to the general public. In addition, digitalisation impacts how we manage projects and also shapes the way we access and interact with information. Organised by the University of Bologna in 2020, the international conference *Italianistica digitale*² was an important opportunity to exchange ideas, presenting and comparing digital resources available for Italian literature, from dictionaries to digital libraries, from monographic portals to critical editions and infra-structures, and also constituting a moment of reflection on the evaluation of these new research products.

The intersection of **digitalisation** and **lexicography** has significantly transformed the world of dictionaries and language resources, a field of intense activity for the Accademia della Crusca, the major Italian institution dedicated to the Italian language, its history, grammar, lexicography, and culture. In this contribution, we delve into the linguistic treasures nowadays also preserved digitally, and no longer only in paper form, by the Accademia,³ including books, manuscripts, and documents (Section 2). We will also point out some projects still in progress (Section 3) and other tools (Section 4).

1. Introduction

In Italy, like in many other countries, the digital revolution has had a great impact on language research as well on the production of dictionaries, terminology databases, and text corpora (see Kirchmeier-Andersen 2011). The introduction of language technology⁴ has helped language institutions in many ways, including

¹ The article was planned by the two authors: Cecilia Robustelli dealt with Sections 1-2.1.3 while Francesca Cialdini covered Sections 2.2-4.3.

² <https://italianisticadigitale-unibo.github.io/>.

³ Founded in 1582-1583 in Florence, the mission of the Accademia della Crusca is to study and preserve the Italian national language. For more information, see the Accademia della Crusca's website at www.accademiadellacrusca.it.

⁴ For a focus on the Italian Language in the Digital Era, see the updated report of the META-NET white paper published in 2012 which was developed as part of the European Language Equality (ELE) project: https://european-language-equality.eu/wp-content/uploads/2022/03/ELE__Deliverable_D1_21_Language_Report_Italian_.pdf.

dealing with the immense number of manuscripts, books, and documents stored in the many libraries scattered around the country which are sometimes difficult to access. At the national level, an important example is the launch, in 1999, of the coordinated digital library project, the *Biblioteca Digitale Italiana* (BDI), in line with similar initiatives carried out in Europe and beyond, by the General Directorate for Libraries, Cultural Institutes and Copyright.⁵ The BDI is now accessible via the portal Internet Culturale,⁶ which also includes – in addition to the National Library Service (SBN – *Servizio bibliotecario nazionale*) – some specialised databases, such as the Census of Manuscripts in Italian Libraries (*Censimento dei manoscritti delle biblioteche italiane* – Manus) and the National Census of 16th Century Italian Books (Edit16).⁷ In Italy, *MediaLibraryOnLine* (MLOL) represents the first Italian network of public, academic, and school libraries for digital lending. To date, there are over 6,500 participating libraries across all Italian regions and 24 foreign countries.

Individual libraries, of course, began the journey towards digitisation in the 1990s. In this paper we focus on the transition to digital by the Accademia della Crusca, a key institution of recognised authority for research on the Italian language, also in relation to its regional varieties, founded in Florence in the second half of the 16th century. Its main early accomplishment was the *Vocabolario degli Accademici della Crusca* (1612), the first dictionary of the Italian language, which was a model for dictionaries of the great European languages. Currently, the Accademia's activities extend beyond historical language knowledge. They actively promote awareness of the ongoing evolution of Italian within the information society (Alisi et al. 2006). Collaborating with both Italian and foreign lexicographic enterprises, the Accademia conducts scientific research through its four research centres: the Centro Studi di Filologia Italiana, publisher of the journal *Studi di Filologia Italiana*;

- a) the Centro Studi di Filologia Italiana, publisher of the journal *Studi di Filologia Italiana*;
- b) the Centro Studi di Lessicografia Italiana, publisher of the journal *Studi di Lessicografia Italiana*;
- c) the Centro Studi di Grammatica Italiana, publisher of the journal *Studi di Grammatica Italiana*;
- d) the Centro di Consulenza Linguistica, which is aimed at institutions, offices, schools, and private citizens through the Accademia's website.⁸

⁵ On the history of the BDI, see Quondam (2021, 137-147).

⁶ www.internetculturale.it/.

⁷ www.iccu.sbn.it/it/internet-cultural/storia-della-biblioteca-digitale-italiana-bdi/index.html; <https://www.internetculturale.it/>.

⁸ See the Statute of the Academy published on the Crusca's website: <https://accademiadella-crusca.it/it/contenuti/statuto-dell'accademia/6956>.

In the 1990s, the Accademia della Crusca initiated its first digital acquisition project: the *Fabbrica dell'Italiano*.⁹ This archive encompasses:

- 1) Dictionaries and grammars: a collection of 2,371 dictionaries and 367 grammars (spanning 1516-2001) preserved in the Biblioteca dell'Accademia della Crusca.
- 2) Manuscripts: 170 manuscripts received by the Crusca for literary competitions in the 19th century.
- 3) Lemmatisation: 9,000 technical terms collected by Cardinal Leopoldo de' Medici.

Subsequently, the Accademia started to collaborate with the Ministero dei Beni e delle Attività Culturali within the national *Biblioteca Italiana Digitale* project (Ragionieri 2015).

Today, the Accademia's library continues to enrich its digital resources, fostering active user participation. In this paper we will describe two main resources of the Accademia:

- a) the digitisation of its book heritage, including manuscripts and incunabula, which are thus made available to the public.
- b) its consulting service aimed at all those seeking grammatical and lexical information as well as clarifications, explanations of linguistic phenomena, and the origin and history of words.

2. The Academy's preservation of digital heritage: *Scaffali Digitali*

Over the years, the Accademia has launched many digitisation projects in order to protect and enhance its library holdings and to facilitate the searchability of materials by the public. In addition, some databases and lexicographic tools useful for studying the evolution of the Italian lexicon from both a synchronic and diachronic perspective have been created and made available.

They can be accessed through the section of the Crusca's website *Scaffali Digitali (Digital Shelves)*,¹⁰ which includes databases and digital archives (Section 2.1), and lexicographic tools (Section 2.2), among which the section *Le Crusche in rete* stands out. We offer some examples of the individual resources accessible from these sections.

⁹ <http://193.205.158.216/fabitaliano2/>.

¹⁰ <https://accademiadellacrusca.it/it/sezioni/scaffale-digitale/25>. A description of *Scaffali digitali* can be found in Biffi/Maraschio (2009) and Maraschio/Marazzini (2021).

2.1 Databases and digital archives

2.1.1 Storia e patrimonio librario dell'Accademia¹¹

This section includes all archives, databases, and repertoires dedicated to the enhancement of the Accademia's book and archival heritage: incunabula, *cinquecentine* (sixteenth century editions), and preparatory materials for the fifth edition of the *Vocabolario degli Accademici della Crusca*.¹²

Of particular interest is the Accademia's incunabula database, containing digital reproductions of the 41 incunabula (approximately 40,000 pages) owned by the Accademia.¹³ Figure 1 illustrates the *Comento sopra la Comedia di Danthe Alighieri poeta fiorentino* by Cristoforo Landino (1481). Each digital volume can be partially “interrogated” because some parts of the text (chapter and paragraph headings) are marked with XML-TEI tags. Users can interact with specific parts of the content, with the digital volumes enhancing accessibility and facilitating targeted exploration.



Fig. 1: Illustration of the first page of *Comento sopra la Comedia di Danthe Alighieri poeta fiorentino*

¹¹ <https://accademiadellacrusca.it/it/contenuti/lista/storia-e-patrimonio-librario-dell-accademia-della-crusca/7491>.

¹² This fifth edition was interrupted in the first half of the 20th century (1923) and has not been continued.

¹³ See <https://incunaboli.accademiadellacrusca.org/index.html>, also for the history of the project.

2.1.2 *Biblioteca digitale*. Fonti normative e descrittive dell'italiano: corpus digitale di testi dal XVI al XIX secolo¹⁴

The *Biblioteca Digitale* is a publicly accessible database containing the complete digitisation of 121 volumes (equivalent to 111,000 images) owned by the Accademia della Crusca, including ancient books relevant to Italian linguistic history such as grammars and dictionaries. In cases where these ancient texts are not held in the Accademia's library, digitisation efforts were coordinated with other libraries, including the National Central Library of Florence and the National Central Library of Rome.

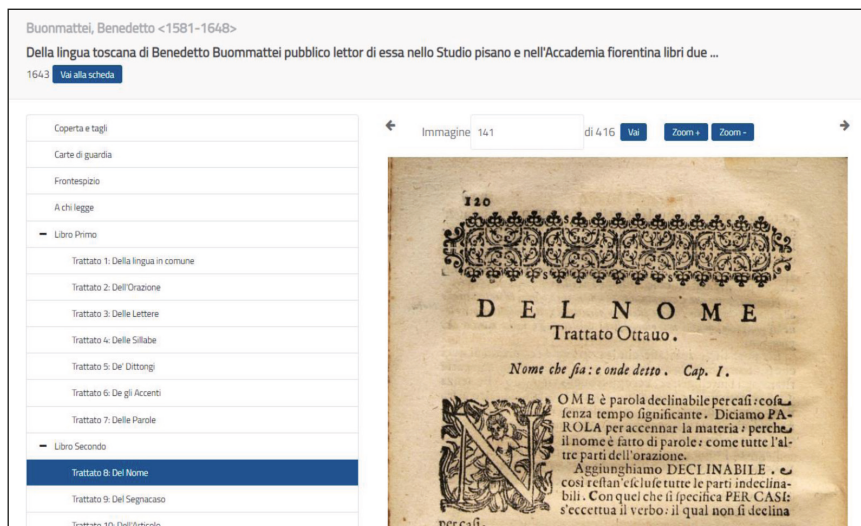


Fig. 2: Illustration from *Della lingua toscana* by Benedetto Buommattei

The digitised editions are organised in four thematic sections:

- 1) Grammars of the Italian language (16th to 19th centuries). Examples include *Prose della volgar lingua* by Pietro Bembo (published in 1525) and *Della lingua toscana* by Benedetto Buommattei (published in 1643) (see Fig. 2).
- 2) Major lexicographic works of the 19th century. This period is often considered the “golden age” of Italian lexicography. Notable works include the *Dizionario della lingua italiana* by Niccolò Tommaseo and Bernardo Bellini (1865-1879) and the *Novo vocabolario della lingua italiana secondo l'uso di Firenze* by Giovan Battista Giorgini and Emilio Broglio (1877-1897), which is based on the linguistic ideas of Alessandro Manzoni.

¹⁴ www.bdcrusca.it/.

- 3) Texts on the role and activities of the Accademia della Crusca (18th century). These texts delve into discussions about the Accademia's functions. For instance, see Paolo Beni's critical work *Anticrusca* (published in 1612).
- 4) Unofficial editions of the *Vocabolario degli Accademici della Crusca*. This section features unofficial versions of the *Vocabolario*, including the edition published in Venice in 1680 and the one by Manuzzi in Florence between 1859 and 1865, which contributed to grammatical and lexical reflection, providing insights into the discussions and reactions surrounding the *Vocabolario* during its publication. Users can navigate digital volumes page by page and perform targeted queries using a detailed index structure.

2.1.3 VIVIT – *Vivi Italiano*¹⁵

VIVIT is an online repository of materials and tools aimed at Italians abroad, particularly those in the second and third generations, representing Italian language and culture. It aims to be a reference point for establishing a strong cultural connection with our country from a distance. The portal provides educational paths and descriptive profiles that highlight the most significant aspects of the Italian language, its history, and its varieties in connection with major historical, artistic, and cultural phenomena. It also offers access to textual databases on contemporary Italian, especially radio and television language, providing a systematic introduction to the Italian spoken today in our country:

- 1) *Lessico di Frequenza dell'Italiano Radiofonico (LIR)*.¹⁶ This database focuses on radio broadcasts in Italian. The initial corpus, created in 1995, comprises 108 hours of content from nine national broadcasters. It includes 64 hours of spoken language (marked according to the XML-TEI standard), with 650,000 occurrences and 86,000 unique forms. A subsequent radio corpus, limited to Rai radio, was established in 2003, containing 32 hours of spoken language, 310,000 occurrences, and 39,000 forms.
- 2) *Lessico dell'Italiano Televisivo (LIT)*.¹⁷ *LIT* is designed for the study of televised Italian. It encompasses 168 hours of content from Rai and Mediaset television broadcasts, captured in 2006 using a statistically representative grid. Researchers can search for specific words or phrases, access quantitative data on frequency, explore contextual usage, and analyse tagged material following the XML-TEI standard (Mauroni/Piotti 2010). Additionally, *LIT* has been enriched with an additional 40 hours of diachronic television speech.

¹⁵ <https://accademiadellacrusca.it/it/contenuti/vivit-vivi-italiano/7467>.

¹⁶ <https://accademiadellacrusca.it/it/contenuti/lessico-di-frequenza-dell-italiano-radiofonico-lir/93>.

¹⁷ <https://accademiadellacrusca.it/it/contenuti/lessico-dell-italiano-televisivo-lit/102>.

- 3) These authentic materials can be consulted by individuals as well as Italian teachers abroad (who often struggle to find systematically accessible and analysable materials for their teaching).
- 3) An electronic dictionary of Italianisms, currently corresponding to the *Dizionario di Italianismi in Francese, Inglese e Tedesco (DIFIT)* compiled by Stammerjohann serves as the initial repository of Italianisms spread abroad. The Academy plans to enrich this repository with contributions from users that will be reviewed by experts in the field. Additionally, a dedicated “cloud” allows community members to share materials (texts, photographs, audiovisual content) on the web.




Fig. 3: VIVIT – Vivi Italiano homepage



Fig. 4: VIVIT – Vivi Italiano sections

Archivi digitali



Da questa sezione del portale si può accedere ad alcuni corpora dell'italiano post-unitario, a strumenti lessicografici e alle numerose banche dati dell'Accademia della Crusca. In particolare **LIT**, **LIR** e **LIS** costituiscono un importante nucleo di base per rappresentare l'italiano scritto e trasmesso. Un **metamatore di ricerca** integra le possibili indagini su questi tre corpora facilitando una prima ricognizione sul lessico italiano. Le banche dati costituiscono un importante strumento per lo studio della lingua italiana, ma allo stesso tempo mettono a disposizione di cultori e studiosi materiale autentico per conoscere meglio la nostra lingua e per insegnarla.

[Interroga il metamatore dei lessici dell'italiano \(LIT-LIR-LIS\)](#)

LIT (Lessico italiano televisivo)
Banca dati interrogabile che raccoglie **168 ore di trasmissioni delle reti RAI e Mediaset** prelevate secondo un campione rappresentativo nel corso del 2006.

LIR (Lessico Italiano Radiofonico)
Banca dati interrogabile che raccoglie **141 ore di trasmissioni di emittenti radiofoniche a diffusione nazionale** prelevate secondo un campione rappresentativo nel 1997 e nel 2003.

LIS (Lessico Italiano Scritto)
Banca dati interrogabile che raccoglie **25 milioni di occorrenze**, distribuite equamente su cinque periodi: 1861-1900, 1901-1922, 1923-1945, 1946-1967, 1968-2001. Rappresenta una versione adattata del **DIACORIS** funzionale all'interrogazione globale di tutte le banche dati presenti nel portale VIVIT.

DIFIT (Dizionario di Italianismi in Francese, Inglese e Tedesco)
Versione elettronica del DIFIT, primo deposito di italianismi incrementabile con i suggerimenti dei consultatori (poi diventato **OIM-Osservatorio degli italianismi del Mondo**).

Banche dati dell'Accademia della Crusca
Un collegamento costante con lo **Scaffale digitale** del sito ufficiale dell'Accademia della Crusca.

Fig. 5: *Archivi digitali*

2.2 Lexicographic tools

2.2.1 Le Crusche in Rete

This section provides access to the five editions of the *Vocabolario degli Accademici della Crusca* spanning from 1612 to 1923. Digitised and meticulously tagged according to XML-TEI standards, this invaluable resource allows users to explore various facets of each entry, including headwords, definitions, literary examples, proverbs, idiomatic expressions, and more. Let us delve into its impressive statistics:

- Total occurrences: over 11 million instances across the editions.
- Forms: more than 200,000 unique word forms.
- Entries: a staggering 142,000 entries.
- Literary examples: enriched with over 400,000 illustrative excerpts.
- Living usage expressions: approximately 6,443 contemporary expressions.
- Proverbs: 2,658 timeless proverbs.
- Locutions: a wealth of 35,779 idiomatic phrases.
- Latin forms: over 32,000 Latin word forms.
- Greek forms: more than 26,000 Greek word forms.

Functioning as a database, the *Vocabolario* unveils the linguistic richness of this lexicographic masterpiece. Notably, the 1612 edition presents a fascinating dual

linguistic layer: ancient Italian (14th-century literary Florentine) and 16th-century Italian (both literary Florentine language and everyday usage). Furthermore, tracing the lexicographic evolution of entries across the five editions provides a captivating lens for studying the history of Italian lexicography.

OPZIONI DI RICERCA

Edizione:	prima ed.(1612)	seconda ed.(1623)	terza ed.(1691)	quarta (1729-1738)	quinta (1863-1923)
Apparati:	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
Dizionario:	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> (solo lemmario)
Giunte:	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	

Ricerca anche su integrazioni crusche

settaggi di ricerca: Standard Considera accenti Considera Minuscole/Maiuscole

Visualizzazione: Voce intera Contesti evidenziati

Ordinamento risultati: Alfabetico Punteggio

libera su Voci

Tipo di ricerca: nei contesti semplice

Contesti

Voce Lemma Definizione Esempio Commento

Cerca nel corpo dei contesti selezionati
 Cerca nei microcontesti

Fig. 6: Search options

2.2.2 The digitisation of the *Grande Dizionario della Lingua Italiana (GDLI)*

This resource stands as a significant project undertaken by the Accademia della Crusca. The digitisation of *GDLI*, a cornerstone of historical dictionaries of the Italian language, offers an invaluable research tool for scholars delving into the intricacies of Italian linguistic history.

Several years ago, the publishing house UTET generously granted the Crusca unrestricted access to the dictionary for scientific purposes. In response, the Accademia embarked on the ambitious task of fully digitising this linguistic treasure trove (Marazzini/Maconi 2018, 100-119) using Optical Character Recognition (OCR), a process that, due to the text’s complexity, encountered challenges related to character recognition and the identification of entries and their components (Sassolini et al. 2021, 160).

To address these intricacies, the Accademia della Crusca initiated a collaborative effort with the Institute of Computational Linguistics “A. Zampolli” at the CNR-Istituto di Linguistica Computazionale “Antonio Zampolli” (CNR-ILC).

Their joint mission involves refining error correction methods and transforming textual content into structured digital data. This endeavour aims to facilitate seamless consultation and integration with other linguistic resources (Sassolini/Biffi 2020, 235-239).

The acquisition of the *GDLI* text represents a pivotal milestone in realising the *Vocabolario Dinamico dell'Italiano Moderno (VoDIM)*, one of the Accademia della Crusca's strategic programmes. In the next section, we will delve deeper into the fascinating world of *VoDIM*.

Currently, the digital *GDLI* offers two search modes beyond the traditional "search by entry" (Fig. 7):

- 1) Search by word form: extracts all contexts where at least one of the entered words appears (Fig. 8).
- 2) Search by sequence: enables users to locate specific portions of text (Fig. 9).

Additionally, users can access the content in both PDF and JPG facsimile format (Fig. 10), enhancing the richness of their exploration.

The screenshot shows the search results for the word "pensiero" in the "Grande dizionario della lingua italiana". The page header includes the UTET logo and the text "Grande dizionario della lingua italiana Prototipo edizione digitale". Below the header, there are navigation links: "PAGINA D'ENTRATA", "RICERCHE", "SALA DI LETTURA", and "GUIDA ALLA CONSULTAZIONE". The search results section shows "Risultati per: pensiero" and "Numero di risultati: 9". Two results are visible: "VOL. XII Pag.1045 - Da PENSIERI a PENSIERO" and "VOL. XII Pag.1046 - Da PENSIERO a PENSIERO". Each result has three icons: a document icon, a magnifying glass icon, and a share icon.

Fig. 7: *Grande Dizionario della Lingua Italiana* search by entry

The screenshot shows the "Ricerca libera" (Free search) interface. It features a search input field with the placeholder text "Cerca...". To the right of the input field is a checkbox labeled "Considera accenti" and a blue button with a magnifying glass icon and the text "Cerca". Below the search field, there are two buttons: "► Elenco forme" and "► Elenco forme per frequenza".

Fig. 8: GDLI search by word form

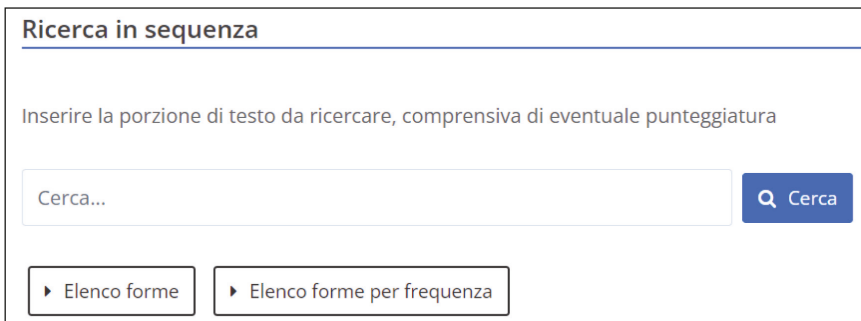


Fig. 9: GDLI search by sequence



Fig. 10: GDLI PDF and JPG format content

2.2.3 ArchiDATA – Archivio di retrodatazioni lessicali¹⁸

ArchiDATA curated by Maconi (2020) serves as a dedicated resource for tracking the earliest documented appearances of words. Focusing primarily on the modern and contemporary lexicon, this database plays a crucial role in refining our understanding of linguistic evolution. As of December 2023, the Lexical Backdating Database had published 11,100 lexical backdatings and 1,670 backdatings of locutions.

Researchers can explore this rich repository through five distinct research modes:

¹⁸ The database is accessible from the website www.archidata.info/.

- 1) Alphabetical search: easily locate specific words.
- 2) Lexical search: investigate word origins and historical contexts.
- 3) Chronological search: trace the chronological emergence of terms.
- 4) Author-based search: explore words associated with renowned figures, including writers and scientists.
- 5) Field of use search: delve into specialised domains where words thrive.

Additionally, users can search for foreign words and regionalisms, enhancing their linguistic exploration (Fig. 11). This dynamic resource contributes significantly to unravelling the intricate tapestry of language origins.

Fig. 11: *ArchiDATA*

3. Work in progress

Three major projects currently still in progress are: the *Vocabolario Dantesco (VD)*, the *Osservatorio degli Italianismi nel mondo (OIM)*, and the *Vocabolario Dinamico dell'Italiano Moderno (VoDIM)*.

3.1 *Vocabolario Dantesco (VD)*

Jointly created by the Accademia della Crusca and the Opera del Vocabolario Italiano (OVI-CNR),¹⁹ the *Vocabolario Dantesco* stands as a tribute to the 750th anniversary of Dante Alighieri's birth. Its purpose is to meticulously compile the lexicon

¹⁹ The *Opera del Vocabolario Italiano (OVI)* is an Institute of the *Consiglio Nazionale delle Ricerche (CNR)* with the institutional task of compiling a historical dictionary of the Italian language. The institute is located in the same premises as the Accademia della Crusca in Florence (www.oivi.cnr.it/en/).

found within Dante's vernacular works, beginning with the monumental *Commedia*. To date, over 1,000 entries have been published.

Methodologically aligned with the lexicographic tools of the OVI, particularly the *Tesoro della Lingua Italiana delle Origini (TLIO)*, the *VD* shares its textual databases and employs the GATTO (Gestione degli Archivi Testuali del Tesoro delle Origini) query software²⁰ (Manni 2020; Manni/Mosti 2022).

Let us explore the structure of a *VD* entry (Fig. 12):

- 1) Headword and grammatical category: each entry begins with the headword, accompanied by its grammatical classification.
- 2) Section details:
 - Index locorum: provides references to specific passages where the word appears in Dante's works.
 - Lexicographic correspondence: links related terms and their meanings.
 - Linguistic note: offers insights into linguistic nuances or historical context.
- 3) Semantic structure of the entry:
 - Definition: systematically records all meanings associated with the word, complete with illustrative examples.
 - Semantic and usage marks: these annotations highlight various components – philosophical, scientific, sector-specific, and more – that constitute Dante's lexicon (Manni/Mosti 2022, 279-281).

The *VD* serves as a linguistic treasure trove, unravelling the intricate layers of Dante's language and enriching our understanding of his literary legacy.

The screenshot shows the interface of the Vocabolario Dantesco. At the top, there is a dark red header with the text 'VOCABOLARIO DANTESCO' in gold and white. To the right of the header are the links 'Home' and 'Il progetto'. Below the header, the entry title 'linguaggio s.m.' is displayed in a large, bold, dark red font. Underneath the title is a navigation bar with icons and labels for 'Frequenza', 'Index locorum', 'Corrispondenze', 'Nota', 'Redattore', and 'Tutto / stampa'. A table shows the frequency of the word in different parts of Dante's works: 'Commedia' with 3 occurrences (3 Inf.) and 'Altre opere' with 1 occurrence (1 Fiore). The main content area contains a definition: 'Lingua parlata da un individuo o da una comunità.' followed by two examples in Italian with their corresponding line numbers and editions: '[1] Inf. 31.78: questi è Nembrotto per lo cui mal coto / pur un linguaggio nel mondo non s'usa.' and '[2] Inf. 31.80: Lasciàno stare e non parliamo a vòto; / ché così è a lui ciascun linguaggio / come l'uso ad altrui, ch'a nullo è noto.'. Below this is a sub-entry '1.1 [Con rif. al fuoco]: crepitio, scoppietto secco e ripetuto (estens.)' with an example: '[1] Inf. 27.14: così, per non aver via né forame / dal principio nel foco, in suo linguaggio / si converbian le parole grame.'

Fig. 12: *Vocabolario Dantesco* entry

²⁰ www.oivi.cnr.it/en/Il-Software.html.

3.2 Osservatorio degli Italianismi nel Mondo (OIM)

The *Osservatorio degli Italianismi nel Mondo (OIM)*, in collaboration with the University of Salzburg, has embarked on a remarkable endeavour to create a comprehensive database encompassing all Italian words and terms of Italian origin that have permeated other languages. This ambitious project celebrates linguistic cross-pollination and sheds light on the global impact of Italian. At its core lies the *Dizionario degli Italianismi in Francese, Inglese e Tedesco (DIFIT)*, a seminal work authored by Stammerjohann and published in 2008. The *DIFIT* serves as the initial nucleus of the *OIM*, meticulously documenting Italianisms in French, English, and German contexts. However, the scope of the project has expanded to include Spanish, Polish, Hungarian, and Portuguese.

An international research consortium is collaborating on this initiative, with units spanning prestigious universities such as Barcelona, Budapest, Dresden, Florence, Krakow, Malta, Milan, Rome, Salzburg, Seville, Toronto, and Warsaw (Heinz 2017).



Fig. 13: *Osservatorio degli Italianismi nel Mondo*

Let us explore the *OIM*'s structure (Fig. 14):

- 1) Word relationships: the *OIM* reconstructs the intricate web of relationships between source Italian words and their manifestations in other languages.
- 2) Database statistics:
 - Total entries: over 12,600 entries currently registered.
 - Categorised entries: approximately 1,500 entries are marked as “complete.”
 - Verification process: entries labelled as “in processing” are awaiting final verification.

- Italianisms: a significant portion – more than 8,900 entries – originates from the *DIFIT*'s collections of Italianisms.
- Constant growth: the database is continually expanding, incorporating entries from diverse linguistic landscapes (Pizzoli/Heinz 2022, 474).

The *OIM* stands as a testament to the enduring influence of Italian across borders, weaving linguistic threads that connect cultures and enrich our shared lexicon.

AFFRESCO s. m.

⊗ Vedi tutti i significati (1)

- [arte/arch./archeol.](#) | Tecnica pittorica che consiste nello stendere i colori su uno strato di intonaco fresco. (GRADIT , 1809 ; La forma oggi non più in uso *a fresco* è attestata dal 1535 (Gradit). L'uso della voce univerbata *affresco* era biasimato dall'Ugolini (1855): "affresco, sostantivamente usato non ha la nostra lingua, ma solo il modo avverbiale a fresco, dipingere a fresco: non dirai dunque affreschi, ma pitture a fresco: potrai però dir freschi per affreschi" (DELLI).)

7

LINGUE IN CUI LA VOCE ITALIANA SI È DIFFUSA

9

FORME DERIVATE DALLA VOCE ITALIANA NELLE ALTRE LINGUE

10

TOTALE SIGNIFICATI ATTESTATI NELLE ALTRE LINGUE

[Versione stampabile](#)

[Legenda](#)

Lingue romanze

Forme totali: 3 Significati totali: 4

>

Francese	Francese	<ul style="list-style-type: none"> • 2 forme in lavorazione • 3 significati 	Vedi le forme
Portoghese	Portoghese Portogallo	<ul style="list-style-type: none"> • 1 forma • 1 significato 	Vedi le forme

Lingue germaniche

Forme totali: 3 Significati totali: 3

>

Inglese	Inglese Gran Bretagna	<ul style="list-style-type: none"> • 2 forme in lavorazione • 2 significati 	Vedi le forme
Tedesco	Tedesco	<ul style="list-style-type: none"> • 1 forma in lavorazione • 1 significato 	Vedi le forme

Altre lingue europee

Forme totali: 2 Significati totali: 2

>

Maltese	Maltese	<ul style="list-style-type: none"> • 1 forma • 1 significato 	Vedi le forme
Ungherese	Ungherese	<ul style="list-style-type: none"> • 1 forma • 1 significato 	Vedi le forme

Fig. 14: *OIM*'s structure

3.3 *Vocabolario Dinamico dell’Italiano Moderno (VoDIM)*

VoDIM is an ongoing project that aims at bringing together Italian texts (written and oral, marked according to the XML-TEI standard) from the Unification of Italy (1861) onwards. The corpus consists of approximately 20 million words concerning several domains: art, songs, law, economy, gastronomy, poetry, and politics as well as journalistic, literary, and scientific prose (Gualdo 2018; Marazzini/Maconi 2018).²¹ A “global search” in the entire corpus is possible, as well as in the thematic corpora.



Fig. 15: *Stazione lessicografica*

The texts of the *VoDIM* corpus are also part of the *Stazione lessicografica*, an integrated reference system that includes dictionaries, databases, and electronic resources.²² The dictionaries include the fifth edition of the *Vocabolario degli Accademici della Crusca* (1863-1923); the *Tommaseo-Bellini* (1865-1879); the *GDLI* (1961-2002); dictionaries of contemporary Italian such as the *Vocabolario Treccani*, the *Sabatini-Coletti* (the edition in the *Corriere della Sera*), the *Nuovo De Mauro*; spelling and pronunciation dictionaries such as the *DOP* (*Dizionario italiano multimediale e multilingue d'ortografia e di pronuncia*) and the *DiPI*

²¹ <https://vodim.accademiadellacrusca.org/>.

²² <https://www.stazionelessicografica.it/>.

(*Dizionario di Pronuncia Italiana*). Also available here are the *Archidata* database with lexical backdating; the *LIR*, *LIT*, and *LIS* databases (*Lessico Italiano Radiofonico*, *Lessico Italiano Televisivo*, and *Lessico Italiano Scritto*); online archives of daily newspapers (*Corriere della Sera* and *Repubblica*); and a database with parliamentary speeches (*Discorsi parlamentari*). To increase the degree to which the corpus was representative, it seemed appropriate to create, in addition, *CoLIWeb* (*Corpus della Lingua Italiana nel Web*). In our pursuit of a more comprehensive and representative corpus, we deemed it fitting to introduce the *CoLIWeb project*, focuses on collecting texts – both written and oral – that span the Italian language after the unification of Italy in 1861 (Fig. 15; Biffi/Ferrari 2020).

4. Other useful tools²³

4.1 *Consulenza linguistica*²⁴

The Accademia della Crusca offers a special service known as the *Consulenza linguistica* (*Language Consulting Service*) to interact with users and clients. While it may not be considered a digital library, it functions as an online tool and has gained popularity among the general public. Established in 2002 and specifically designed for those seeking grammatical and lexical clarification, it provides explanations on linguistic facts, delves into the origin and history of words, and addresses various language-related inquiries. Common areas of concern include morphosyntax (especially the usage of specific verbs and matters of agreement), the lexicon (including dialect forms, regionalisms, special languages, and foreign terms), and graphic-phonetic issues (D’Achille 2022; D’Achille/Biffi 2022, 18). This digital platform complements the traditional service offered for years in the pages of the periodical *La Crusca per Voi* founded by the linguist Giovanni Nencioni in 1990.

Consisting of linguists, the editorial team carefully reviews incoming questions and selects those with recurring themes or of widespread interest. Over the past decade, tens of thousands of inquiries have been received and many have received personalised responses. The *Language Consulting Service* section of the website currently boasts 1,358 answers, which can be conveniently searched using tags and keywords (Fig. 16).

²³ All three services (3.1-3) can be accessed by clicking on *Lingua Italiana* from the main menu on the homepage of the Accademia.

²⁴ <https://accademiadellacrusca.it/it/contenuti/consulenza-linguistica/6945>.

Fig. 16: *Consulenza linguistica*

4.2 *Stazione Bibliografica*²⁵

The aim of the section *Stazione Bibliografica (Bibliographic Station)* is to provide a bibliographical overview of Italian linguistics at various levels. It is organised into two parts:

- *Bibliografia Essenziale (Essential Bibliography)*: a compilation of fundamental studies on Italian linguistics.
- *Novità Bibliografiche (Recent Publications)*: here, the latest publications related to linguistic topics are highlighted.

4.3 *Parole Nuove*²⁶

The *Parole Nuove (New Words)* section is a fascinating collection of words that have not yet been included in dictionaries. These words are carefully selected by the editorial staff, without any regulatory intent. The selection process involves monitoring mass media sources (such as national newspapers) and considering user reports. *Parole Nuove* serves as a valuable tool for understanding new words that are circulating in both spoken and written language.

Figure 17 lists some of the intriguing words that gained prominence and were studied in 2023:

²⁵ <https://accademiadellacrusca.it/contenuti/stazione-bibliografica/6944>.

²⁶ <https://accademiadellacrusca.it/contenuti/parole-nuove/7092>.

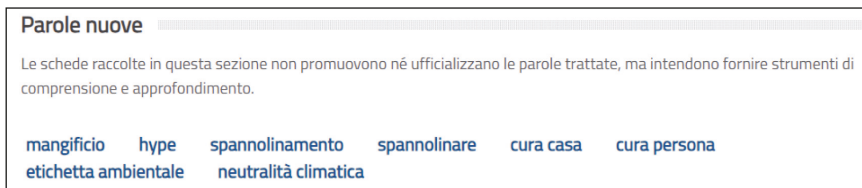


Fig. 17: Some intriguing words in *Parole nuove*

4.4 *Italiano Digitale – La rivista delle Crusca in rete*²⁷

Lastly, we would like to recall the magazine *Italiano Digitale* in electronic format, which, in addition to hosting studies and essays, collects the best publications that have appeared on the Accademia della Crusca's website: linguistic advice sheets provided to readers, lexicographic sheets dedicated to the most recent Italian words, discussion topics proposed by academics, and articles of linguistic and historical-linguistic interest that the Accademia chooses to highlight. The aim of *Italiano digitale* is to provide all those interested with access to the Accademia della Crusca's scientific contributions in PDF format. Each issue is accompanied by a summary of the most significant activities and initiatives involving the Accademia and its members.

References

- Alisi, T. et al. (2006): Advanced search facilities for accessing Crusca Academy of Italian Language. In: Cappellini, V./Hemsley, J. (eds.): *Electronic imaging & the Visual Arts EVA 2006 Florence Proceedings*. Bologna: Pitagora Editrice, 164-169.
- Biffi, M./Fanfani, M. (2006): La Lessicografia della Crusca in Rete. In: Corino, E./Marello C./Onesti C. (eds.), *Proceedings of XII Euralex International Congress, Torino, September 2006*. Alessandria: Dell'Orso, 409-416.
- Biffi, M./Ferrari, A. (2020): Progettare e ideare un corpus dell'italiano nella rete: il caso del CoLIWeb. In: *Studi di Lessicografia Italiana*, XXXVII, 357-374.
- Biffi, M./Maraschio, N. (2009): Strumenti digitali dell'Accademia della Crusca. In: Magherini, S. (ed.), *Tradizione e Modernità. Archivi digitali e strumenti di ricerca*. Florence: Società Editrice Fiorentina, 115-146.
- D'Achille, P. (2022): Aspetti e problemi del lessico italiano nel Servizio di Consulenza linguistica dell'Accademia della Crusca. In: González Royo, C./Nappi, P. (eds.): *Parole a confronto. Lessicografia, traduzione e didattica tra italiano e spagnolo*. Frankfurt a.M.: Peter Lang, 49-67.

²⁷ <https://id.accademiadellacrusca.org/>.

- D'Achille, P./Biffi, M. (2022): Premessa. In: Accademia della Crusca (eds.): *Giusto, sbagliato, dipende. Le risposte ai tuoi dubbi sulla lingua italiana*. Milan: Mondadori, 17-21.
- Gualdo, R. (2018): Un nuovo Vocabolario dinamico dell'italiano. Il lessico specialistico e settoriale. In: *Studi di lessicografia italiana XXXV*, 193-216.
- Heinz, M. (ed.) (2017): *Osservatorio degli italianismi nel mondo. Punti di partenza e nuovi orizzonti*. Florence: Accademia della Crusca.
- Kirchmeier-Andersen, S. (2011): Language technology for language institutions. What kind of technology do languages institutions use, what kind of resources can they provide? In: Stickel, G./Varadi, T. (eds.): *Language, languages and new technologies: ICT in the service of languages*. (= Duisburger Arbeiten zur Sprach- und Kulturwissenschaft 87). Frankfurt a.M.: Peter Lang, 21-32.
- Maconi, L. (ed.) (2020): *Laboratorio di ArchiDATA 2020. Retrodatazioni lessicali: storia di cose e di parole*. Florence: Accademia della Crusca.
- Manni, P. (ed.) (2020): «S'i' ho ben la parola tua intesa». *Atti della giornata di presentazione del Vocabolario Dantesco (Firenze, Accademia della Crusca, 1° ottobre 2018)*. Florence: Accademia della Crusca.
- Manni, P./Mosti, R. (2022): Per Dante. Il VD e i corpora dell'italiano antico. In: Cresti, E./Moneglia M. (eds.): «Corpora e Studi Linguistici». *Atti del LIV Congresso della Società di Linguistica Italiana (Online, 8-10 settembre 2021)*. Milan: Officinaventuno, 275-293.
- Maraschio, N./Marazzini, C. (2021): Gli scaffali digitali della Crusca. In: *Italianistica digitale*, 20, 2, 91-101.
- Marazzini, C./Maconi, L. (2018): Il 'Vocabolario dinamico dell'italiano moderno' rispetto ai linguaggi settoriali. Proposta di voce lessicografica per il redigendo VoDIM. In: *Italiano digitale VII*, 4, 100-119.
- Mauroni, E./Piotti, M. (eds.) (2010): *L'italiano televisivo 1976-2006*. Florence: Accademia della Crusca.
- Pizzoli, L./Heinz, M. (2022): Il progetto OIM (Osservatorio degli Italianismi nel Mondo). In: *Italiano LinguaDue* 14, 2, 471-487.
- Quondam, Amedeo (2021): Memorie per una storia dell'italianistica digitale: Biblioteca italiana. In: *Griseldaonline* 20, 138-147.
- Ragionieri, D. (2015): *La Biblioteca dell'Accademia della Crusca*. Florence: Accademia della Crusca e Manziiana/Vecchiarelli Editore Srl.
- Robustelli, C. (2013): Il Vocabolario dell'Accademia della Crusca e i primi grandi vocabolari delle lingue europee. In: Stickel, G./Varadi, T. (eds.): *Lexical Challenges in a Multilingual Europe. Contributions to the Annual Conference 2012 of EFNIL in Budapest*. Frankfurt a.M.: Peter Lang, 127-137.

- Sassolini, E. et al. (2021): La digitalizzazione del GDLI: un approccio linguistico per la corretta acquisizione del testo?. In: Boschetti, F./Del Grosso, A.M./Salvatori, E. (eds.): *AIUCD 2021-DH per la società: e-guaglianza, partecipazione, diritti e valori nell'era digitale*. AIUCD Associazione per l'Informatica Umanistica e la Cultura Digitale, 159-166.
- Sassolini, E./Biffi, M. (2020): Strategie e metodi per il recupero di dizionari storici. In: *IX Convegno annuale AIUCD. La svolta inevitabile: sfide e prospettive per l'informatica umanistica*. Milan: Università Cattolica del Sacro Cuore, 235-239.
- Sessa, M. (1980): Il “rovesciamento” del ‘Vocabolario della Crusca. In: *Bollettino d'informazioni del Centro di elaborazione automatica di dati e documenti storico-artistici della Scuola Normale Superiore di Pisa* I, I, 41-51.
- Stammerjohann, H. (2008): *Dizionario di italianismi in francese, inglese e tedesco*. Florence: Accademia della Crusca.

