

# Ridire 2.0

## La nuova *release* del web Corpus promosso dalla SILFI

RIDIRE.it (Risorsa Dinamica Italiana di Rete) è un web corpus che documenta l'italiano utilizzato in rete nei domini più rappresentativi della lingua e della cultura italiana, frutto di un progetto promosso dalla SILFI, in collaborazione con un consorzio di università italiane,<sup>1</sup> e finanziato dal Fondo Italiano per la Ricerca di Base (FIRB).

RIDIRE è stato rilasciato per la prima volta dieci anni fa, in un'infrastruttura di rete dotata di funzioni computazionali di ricerca proprie della linguistica dei corpora necessarie a evidenziare le particolarità dell'uso linguistico e della fraseologia italiana, sia a fini di ricerca e come strumento didattico sia per il consolidamento del possesso dell'italiano da parte degli apprendenti.

Ridire 2.0 ripropone il web corpus in una nuova *release*, in cui il rumore causato dall'estrazione automatica dei testi dal web (duplicazioni dei dati, problemi di formato, *boilerplate*) è stato drasticamente ridotto e l'infrastruttura computazionale di ricerca, aggiornata agli standard contemporanei, consente oggi un accesso ai dati più semplice e veloce.

I documenti raccolti in RIDIRE costituiscono una fonte rilevante per lo studio della fraseologia italiana, sia generale sia nei suoi ambiti settoriali. La base dati interrogabile ha una dimensione di circa 1.3 miliardi di parole e è strutturata in domini divisi in due tipologie

<b>Domini Semantici:</b> i campi dell'eccellenza italiana in cui le scelte linguistiche sono dettate dal contenuto	<b>Domini Funzionali:</b> contesti d'uso in cui la lingua ha convenzioni legate alla funzione sociale
<ul style="list-style-type: none"><li>• Letteratura</li><li>• Moda</li><li>• Design-architettura</li><li>• Cucina</li><li>• Sport</li><li>• Religione</li><li>• Arti figurative</li><li>• Cinema</li><li>• Musica</li></ul>	<ul style="list-style-type: none"><li>• Informazione</li><li>• Economia e affari</li><li>• Amministrazione e legislazione</li></ul>

### Il design di RIDIRE

Ridire 2.0 è mantenuto in rete dal Laboratorio di Informatica Umanistica del Dipartimento di Lettere e Filosofia dell'Università di Firenze (LIU) e è liberamente disponibile all'indirizzo

<http://corpora.dilef.unifi.it/query?corpname=ridire>

<sup>1</sup> Università di Firenze (LABLITA e Dipartimento Sistemi e Informatica), l'Università di Torino (Dipartimento Scienze Letterarie e Filologiche), l'Università di Roma3 (Dipartimento di Italianistica), l'Università di Napoli (Dipartimento di Filologia Moderna) e l'Università di Siena (Dipartimento di Economia).

## Le funzioni di ricerca di Ridire 2.0

La disponibilità di strumenti computazionali per l'estrazione di informazioni vive da grandi corpora rappresentativi dell'uso è un'importante opportunità di conoscenza e costituisce il portato naturale dell'evoluzione tecnologica.

Le occorrenze del corpus RIDIRE sono lemmatizzate e annotate per Parte del Discorso (PoS). Le funzioni di ricerca sono implementate nell'infrastruttura KonText<sup>2</sup>, che permette di effettuare ricerche sull'intero corpus o sui suoi sotto-corpora, estraendo Liste di frequenza, Concordanze, Collocazioni e *Colligation*.

The screenshot displays the KonText search interface. At the top, there are tabs for 'Advanced query', 'Keyboard', and 'Query interpretation'. Below the search input field, a tip explains that a color highlighted token with a question mark symbol indicates an additional specification available. The interface includes several sections for configuring the search:

- Specify parameters:** Includes 'Match case', 'Allow regular expressions', and 'Default attribute: word'.
- Specify context:** A section for defining the search context.
- Restrict search:** A section for restricting the search to specific corpora.

Under the 'Restrict search' section, there is a 'Save as a subcorpus draft' button and four panels for selecting corpora:

- doc.id:** Includes an 'Exact value...' input field.
- doc.jobname:** Includes an 'Exact value...' input field.
- doc.source:** A list of source types with their frequencies:

<input type="checkbox"/> ===NONE===	14,162,245
<input type="checkbox"/> doc	17,437,727
<input type="checkbox"/> html	914,108,890
<input type="checkbox"/> pdf	429,285,016
<input type="checkbox"/> rtf	326,250
- doc.domain:** A list of domains with their frequencies:

<input type="checkbox"/> ===NONE===	14,162,245
<input type="checkbox"/> Amministrazione e Legislazione	373,688,122
<input type="checkbox"/> Architettura e Design	90,161,631
<input type="checkbox"/> Arti figurative	71,817,586
<input type="checkbox"/> Cucina	28,494,278
<input type="checkbox"/> Economia e Affari	137,394,729
<input type="checkbox"/> Informazione	177,991,674
<input type="checkbox"/> Letteratura e Teatro	127,482,254
<input type="checkbox"/> Moda	19,539,320
<input type="checkbox"/> Musica	65,590,400
<input type="checkbox"/> Non classificato	30,840,021
<input type="checkbox"/> Religione	92,854,041
<input type="checkbox"/> Sport	116,808,900

### La schermata di KonText per la scelta dei corpora su cui effettuare le ricerche

Una videoguia all'utilizzo delle funzioni di ricerca implementate in KonTexte realizzata da LIU e è disponibile qui [http://lablita.eu/downloads/userguide\\_ridire2.0.mp4](http://lablita.eu/downloads/userguide_ridire2.0.mp4)

<sup>2</sup> Machálek, T. (2020). KonText: Advanced and flexible corpus query interface. In *Proceedings of the Twelfth Language Resources and Evaluation Conference* (pp. 7003-7008).