

Word Integration

Amedeo Leo¹ Alessio Petrozziello¹
Simone Romano¹

¹Università degli studi di Salerno

Presentazione progetto Gestione Avanzata dei Dati
2014/2015

- 1** Introduzione
- 2** Fonti
- 3** Wrapper
- 4** Schema globale
- 5** Tecnologie

- 1** Introduzione
- 2 Fonti
- 3 Wrapper
- 4 Schema globale
- 5 Tecnologie

Obiettivo:

- Generare informazioni di diverso tipo su una parola (italiana/inglese) data in input.



The screenshot shows the main interface of the 'Word Integration' application. At the top, there is a dark navigation bar with the text 'Word Integration', 'Home', and 'About'. Below this, the title 'Word Integration' is displayed in a large font, with a blue button labeled 'Learn more ...' underneath. A search section follows, starting with the label 'Parola/Word:' and a text input field containing the placeholder 'Inserisci la parola/Insert the word...'. Below the input field are two radio buttons for language selection: 'Italiano' (selected) and 'Inglese'. Underneath are several checkboxes for search options: 'Sinonimi', 'Contrari', 'Rime', 'Contesti', 'Wikipedia', 'Dialetti', 'Traduzione', and 'Modi di dire'. A blue 'Invia' button is positioned below these options. At the bottom of the page, there is a light gray footer containing the names 'Amedeo Leo', 'Alessio Petrozziello', and 'Simone Romano', each accompanied by a black and white LinkedIn logo icon.

Word Integration [Home](#) [About](#)

Traduzioni


cane (da en a it)


Wikipedia




The domestic dog (*Canis lupus familiaris*) is a canid that is known as man's best friend. The dog was the first domesticated animal and has been widely kept as a working, hunting, and pet companion. According to recent coarse estimates, there are currently between 700 million and one billion dogs, making them the most abundant member of order Carnivora in the world.

Autori

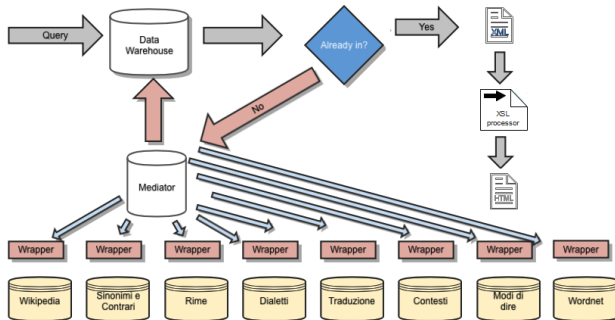
Amedeo Leo, Alessio Petrozziello e Simone Romano sono studenti dell'Università degli Studi di Salerno. L'applicazione web Word Integration è stata creata come progetto per il corso di Gestione Avanzata dei Dati, tenuto dal Professor Gennaro Costagliola, afferente al Dipartimento di Informatica della stessa Università.

Il codice sorgente si può scaricare qui: 

La documentazione è visibile qui: 

Amedeo Leo  Alessio Petrozziello  Simone Romano 

Architettura ibrida



- 1 Introduzione
- 2 Fonti**
- 3 Wrapper
- 4 Schema globale
- 5 Tecnologie

Possiamo raggruppare le fonti utilizzate in 2 categorie:

1 Italiane

- <http://www.sinonimi-contrari.it> STATICA
- <http://www.cercarime.it/?rima> STATICA
- <http://www.dialettando.com/dizionario/> CAMBI PERIODICI
- Google translator API STATICA
- Wikipedia API STATICA
- Modi di dire STATICA
- <http://web-ngram.research.microsoft.com/similarity/> CAMBI PERIODICI

2 Inglese

- Google translator API STATICA
- WordNet API STATICA
- Wikipedia API STATICA

- SinonimiContrari(Parola, Sinonimo, Contrario)
- Cercarime(Parola, Rima)
- Contesti(Parola, Contesto)
- Dialetti(Parola, Dialetto, Regione)
- GoogleTranslator(Parola, Da, A, Traduzione)
- WordNet(Parola, Overview, Synonyms)
- Wikipedia(Parola, Descrizione)
- Modi(Parola, Modo)

- 1 Introduzione
- 2 Fonti
- 3 Wrapper**
- 4 Schema globale
- 5 Tecnologie

Come si evince dall'architettura mostrata, a runtime avvengono i seguenti passaggi:

- Ricerca dei dati nel Database:
 - **Dati presenti:** Generazione output
 - **Dati non presenti:** scraping che mira a popolare il database e generazione output



Strumenti usati

Gli strumenti usati per ricavare le *xpath* sono i seguenti:

- **Curl** per:
 - Estrazione html
 - Generazione xhtml
- **Pentaho Data Integration (Kettle)** per ottenere le *xpath*
- **Exist** per testare le *xpath*

Vediamo dunque le *xpath* di scraping per ciascuna fonte.



Sinonimi e contrari

Per la fonte *<http://www.sinonimi-contrari.it>* le xpath utilizzate sono:

- `/*[name()='html']/*[name()='body']/*[name()='div']/*[name()='div']/*[name()='div']
/*[name()='div'][1]/*[name()='ol']/*[name()='li']/*[name()='span']/*[name()='a']/text()`
- `/*[name()='html']/*[name()='body']/*[name()='div']/*[name()='div']/*[name()='div']
/*[name()='div'][3]/*[name()='ol']/*[name()='li']/*[name()='span']/*[name()='a']/text()`



SINONIMI - CONTRARI

Sinonimi e contrari della lingua Italiana

Per la fonte <http://www.cercarime.it/?rime> le xpath utilizzate sono:

- VECCHIO
`//*[@name()='html']/*[name()='body']/*[name()='div']/*[name()='div']/*[name()='div']/*[name()='div']
/*[name()='ul']/*[name()='li']/*[name()='p']/*[name()='i']/*[name()='a']`
- NUOVO `//*[@name()='div']/*[name()='div'][1]/*[name()='ul'][1]/*[name()='li']/*[name()='p']/*[name()='i']
/*[name()='a']/text()`

Cerca Rime

Inserisci la rima che vuoi cercare:

Cerca Rima

Per la fonte

<http://web-ngram.research.microsoft.com/similarity/> è stato sufficiente utilizzare la seguente sintassi *php*:

- *valori* = *explode*("",a);
 valori[0] = *substr*(*valori*[0], 3);

Clustering words based on context similarity

Per la fonte <http://www.dialettando.com/dizionario/> è stata utilizzata la seguente espressione xpath:

- ```
//*[name()='html']/*[name()='body']/*[name()='center']/*[name()='table']/*[name()='tr'][5]
/*[name()='td'][2]/*[name()='table']/*[name()='tr'][1]/*[name()='td'][2]/*[name()='center']/*[name()='table']
/*[name()='tr']/*[name()='td']/text()
```



## Google translate

Per la traduzione sono state utilizzate delle API che interagiscono con Google translate:

- *obj* = *newGoogleTranslate(\$da, \$a);*  
*obj->translate(\$parola);*



Per la fonte *wordnet.princeton.edu/* sono state utilizzate delle API:

- `"/usr/local/bin/wn ". $parola ." -over > scraping.html";`  
`"/usr/local/bin/wn ". $parola ." -synsn > scraping.html";`

## WordNet Search - 3.1

- [WordNet home page](#) - [Glossary](#) - [Help](#)

Word to search for:

## Wikipedia

Per la fonte *www.wikipedia.com/* sono state utilizzate delle API che estraggono il contenuto dell'output di wikipedia e lo ripuliscono restituendo una stringa:

- *html = file\_get\_contents(url);*  
*info = get\_text(html);*  
*info = strip\_tags(info);*  
*info = clean\_input(info);*



**WIKIPEDIA**  
L'enciclopedia libera

### Benvenuti su Wikipedia

L'enciclopedia libera e collaborativa

[Sfoggia l'indice](#) [> Consulta il sommario](#) [🌐 Naviga tra i portali tematici](#)

## Modi di dire

Per la fonte

*<http://dizionari.corriere.it/dizionario-modi-di-dire/>* è stata  
utilizzata la seguente xpath:

- `/*[name()='html']/*[name()='body']/*[name()='div'][9]/*[name()='div']/*[name()='div']  
/*[name()='div'][2]/*[name()='div'][1]/*[name()='div'][2]/*[name()='div']/*[name()='a']/text()`

### Dizionario dei Modi di Dire



Edizione online tratta da:

***Dizionario dei Modi di Dire*** HOEPLI Editore  
della Lingua Italiana

CERCA

- 1 Introduzione
- 2 Fonti
- 3 Wrapper
- 4 Schema globale**
- 5 Tecnologie

## Schema globale

Lo schema globale del sistema di interrogazione può essere riscritto come segue:

- SchemaGlobale(Parola, Sinonimo, Contrario, Rima, Dialetto, Regione, Traduzione, Da, A, Overview, Synonyms, Descrizione, Wikipedia, Modo) :-

SinonimiContrari(Parola, Sinonimo, Contrario) ∨ Cercarime(Parola, Rima) ∨ Contesti(Parola, Contesto) ∨ Dialetti(Parola, Dialetto, Regione) ∨ GoogleTranslator(Parola, Da, A, Traduzione) ∨ WordNet(Parola, Overview, Synonyms) ∨ Wikipedia(Parola, Descrizione) ∨ Modi(Parola, Modo)



Riportiamo un **sottoinsieme di query** effettuabili sul sistema:

- 1 Ricerca di un sinonimo e di una rima della parola

"gatto"

$q(\text{Sinonimo, Rima}) = \text{SchemaGlobale}(\text{Parola, Sinonimo, Contrario, Rima, Dialecto, Regione, Traduzione, Da, A, Overview, Synonyms, Descrizione, Wikipedia, Modo}) \wedge \text{Parola} = \text{'gatto'} \wedge \text{Da} = \text{'it'} \wedge \text{A} = \text{'en'}$

- 2 Ricerca tutte le informazioni della parola italiana

"bello"

$q(\text{Sinonimo, Contrario, Rima, Dialecto, Modo, Wikipedia}) = \text{SchemaGlobale}(\text{Parola, Sinonimo, Contrario, Rima, Dialecto, Regione, Traduzione, Da, A, Overview, Synonyms, Descrizione, Wikipedia, Modo}) \wedge \text{Parola} = \text{'gatto'} \wedge \text{Da} = \text{'it'} \wedge \text{A} = \text{'en'}$

- 3 Ricerca informazioni su Informazioni Wikipedia e la traduzione della parola "dog"

$q(\text{Wikipedia, Traduzione}) = \text{SchemaGlobale}(\text{Parola, Sinonimo, Contrario, Rima, Dialecto, Regione, Traduzione, Da, A, Overview, Synonyms, Descrizione, Wikipedia, Modo}) \wedge \text{Parola} = \text{'dog'} \wedge \text{Da} = \text{'en'} \wedge \text{A} = \text{'it'}$

- 4 Prendere i dialetti e le rime della parola "cavallo"

$q(\text{Dialecto, Rima}) = \text{SchemaGlobale}(\text{Parola, Sinonimo, Contrario, Rima, Dialecto, Regione, Traduzione, Da, A, Overview, Synonyms, Descrizione, Wikipedia, Modo}) \wedge \text{Parola} = \text{'cavallo'} \wedge \text{Da} = \text{'it'} \wedge \text{A} = \text{'en'}$





## Di seguito la riformulazione GAV delle interrogazioni:

- 1  $q'(\text{Sinonimo}, \text{Rima}) = \text{SinonimiContrari}(\text{Parola}, \text{Sinonimo}, \text{Contrario}) \vee \text{CercaRime}(\text{Parola}, \text{Rima}) \wedge \text{Parola} = \text{"gatto"} \wedge \text{Da} = \text{"it"} \wedge \text{A} = \text{"en"}$
- 2  $q'(\text{Sinonimo}, \text{Contrario}, \text{Rima}, \text{Dialecto}, \text{Regione}, \text{Modo}, \text{Wikipedia}) = \text{SinonimiContrari}(\text{Parola}, \text{Sinonimo}, \text{Contrario}) \vee \text{CercaRima}(\text{Parola}, \text{Rima}) \vee \text{Dialetti}(\text{Parola}, \text{Dialecto}, \text{Regione}) \vee \text{Modi}(\text{Parola}, \text{Modo}) \vee \text{Wikipedia}(\text{Parola}, \text{Descrizione}) \wedge \text{Parola} = \text{"bello"} \wedge \text{Da} = \text{"it"} \wedge \text{A} = \text{"en"}$
- 3  $q'(\text{Wikipedia}, \text{Trauzione}) = \text{Wikipedia}(\text{Parola}, \text{Descrizione}) \vee \text{GoogleTranslator}(\text{Parola}, \text{Da}, \text{A}) \wedge \text{Parola} = \text{"dog"} \wedge \text{Da} = \text{"en"} \wedge \text{A} = \text{"it"}$
- 4  $q'(\text{Dialecto}, \text{Rima}) = \text{CercaRime}(\text{Parola}, \text{Rima}) \vee \text{Dialetti}(\text{Parola}, \text{Dialecto}, \text{Regione}) \wedge \text{Parola} = \text{"cavallo"} \wedge \text{Da} = \text{"it"} \wedge \text{A} = \text{"en"}$

- **SinonimiContrari(Parola, Sinonimo, Contrario):-**  
SchemaGlobale(Parola, Sinonimo, Contrario, Rima, Dialetto, Regione, Traduzione, Da, A, Overview, Synonyms, Descrizione, Wikipedia, Modo)
- **Cercarime(Parola, Rima):-** SchemaGlobale(Parola, Sinonimo, Contrario, Rima, Dialetto, Regione, Traduzione, Da, A, Overview, Synonyms, Descrizione, Wikipedia, Modo)
- **Contesti(Parola, Contesto):-** SchemaGlobale(Parola, Sinonimo, Contrario, Rima, Dialetto, Regione, Traduzione, Da, A, Overview, Synonyms, Descrizione, Wikipedia, Modo)
- **Dialetti(Parola, Dialetto, Regione):-** SchemaGlobale(Parola, Sinonimo, Contrario, Rima, Dialetto, Regione, Traduzione, Da, A, Overview, Synonyms, Descrizione, Wikipedia, Modo)
- **GoogleTranslator(Parola, Da, A, Traduzione):-**  
SchemaGlobale(Parola, Sinonimo, Contrario, Rima, Dialetto, Regione, Traduzione, Da, A, Overview, Synonyms, Descrizione, Wikipedia, Modo)
- **WordNet(Parola, Overview, Synonyms):-** SchemaGlobale(Parola, Sinonimo, Contrario, Rima, Dialetto, Regione, Traduzione, Da, A, Overview, Synonyms, Descrizione, Wikipedia, Modo)
- **Wikipedia(Parola, Descrizione):-** SchemaGlobale(Parola, Sinonimo, Contrario, Rima, Dialetto, Regione, Traduzione, Da, A, Overview, Synonyms, Descrizione, Wikipedia, Modo)
- **Modi(Parola, Modo):-** SchemaGlobale(Parola, Sinonimo, Contrario, Rima, Dialetto, Regione, Traduzione, Da, A, Overview, Synonyms, Descrizione, Wikipedia, Modo)

- 1  $q(\text{Sinonimo, Rima}) :- \text{SchemaGlobale}(\text{Parola, Sinonimo, Contrario, Rima, Dialetto, Regione, Traduzione, Da, A, Overview, Synonyms, Descrizione, Wikipedia, Modo}) \wedge \text{Parola} = \text{"gatto"} \wedge \text{Da} = \text{"it"} \wedge \text{A} = \text{"en"}$
- 2  $q(\text{Sinonimo, Contrario, Rima, Dialetto, Regione, Modo, Wikipedia}) :- \text{SchemaGlobale}(\text{Parola, Sinonimo, Contrario, Rima, Dialetto, Regione, Traduzione, Da, A, Overview, Synonyms, Descrizione, Wikipedia, Modo}) \wedge \text{Parola} = \text{"bello"} \wedge \text{Da} = \text{"it"} \wedge \text{A} = \text{"en"}$
- 3  $q(\text{Wikipedia, Traduzione}) :- \text{SchemaGlobale}(\text{Parola, Sinonimo, Contrario, Rima, Dialetto, Regione, Traduzione, Da, A, Overview, Synonyms, Descrizione, Wikipedia, Modo}) \wedge \text{Parola} = \text{"dog"} \wedge \text{Da} = \text{"en"} \wedge \text{A} = \text{"it"}$
- 4  $q(\text{Dialetto, Rima}) :- \text{SchemaGlobale}(\text{Parola, Sinonimo, Contrario, Rima, Dialetto, Regione, Traduzione, Da, A, Overview, Synonyms, Descrizione, Wikipedia, Modo}) \wedge \text{Parola} = \text{"cavallo"} \wedge \text{Da} = \text{"it"} \wedge \text{A} = \text{"en"}$

## Bucket example

- $q(\text{Sinonimo, Rima}) :- \text{SchemaGlobale}(\text{Parola, Sinonimo, Contrario, Rima, Dialetto, Regione, Traduzione, Da, A, Overview, Synonyms, Descrizione, Wikipedia, Modo}) \wedge \text{Parola} = \text{"gatto"} \wedge \text{Da} = \text{"it"} \wedge \text{A} = \text{"en"}$   
**Bucket**=SchemaGlobale{V1('gatto', Sinonimo, Contrario), V2('gatto', Rima), V3('gatto', Contesto), V4('gatto', Dialetto, Regione), V5('gatto', Da, A, Traduzione), V6('gatto', Overview, Synonyms), V7('gatto', Descrizione), V8('gatto', Modo)}

- 1 Introduzione
- 2 Fonti
- 3 Wrapper
- 4 Schema globale
- 5 Tecnologie**

## Mapping con il codice e tecnologie utilizzate

Per realizzare l'applicazione web descritte abbiamo utilizzato:

- **php** come linguaggio di programmazione
- database **sql** per memorizzare le informazioni informazioni statiche
- **xml**, **xpath** ed **xslt** per estrarre i dati in real time dal database e mostrarli a video



Word  
Integration

Leo,  
Petrozziello,  
Romano

Introduzione

Fonti

Wrapper

Schema  
globale

Tecnologie

**End**

**Grazie per l'attenzione**