

# Dizionario Ontologico delle Espressioni Polirematiche dell'Italiano

Ph.D student Michele Stingo
DIPARTIMENTO DI SCIENZE POLITICHE E DELLA COMUNICAZIONE
UNIVERSITÀ DEGLI STUDI DI SALERNO - NETWORK CONTACTS
SRL

#### Obiettivi realizzativi

Creazione di una risorsa lessicografica digitalizzata per le MWE

Raccolta ed integrazione di diverse fonti informative

Paradigma Open Data



#### Espressioni Polirematiche



Creazione di una risorsa lessicografica digitalizzata per le MWE

Raccolta ed integrazione di diverse fonti informative

Paradigma Open Data

"gruppo di parole che ha un significato unitario, non desumibile da quello delle parole che lo compongono, sia nell'uso corrente sia nei linguaggi tecnico-specialistici" (De Mauro, 2000)

Ferro da stiro

Conflitto d'interesse

#### Dizionario Ontologico



Creazione di una risorsa lessicografica digitalizzata per le MWE

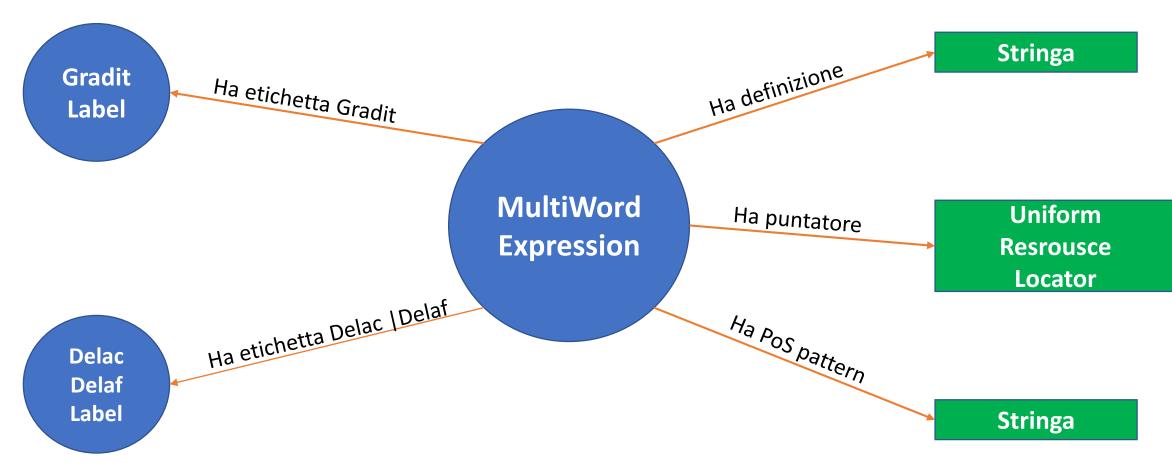
Raccolta ed integrazione di diverse fonti informative

Paradigma Open Data

Un'ontologia è una descrizione formale ed esplicita dei concetti (classi) appartenenti ad un dominio della conoscenza, delle caratteristiche (proprietà) attribuibili ad ogni concetto e delle relative restrizioni. Quando un'ontologia è popolata dai dati, essa assume anche la funzione di Knowledge Base.

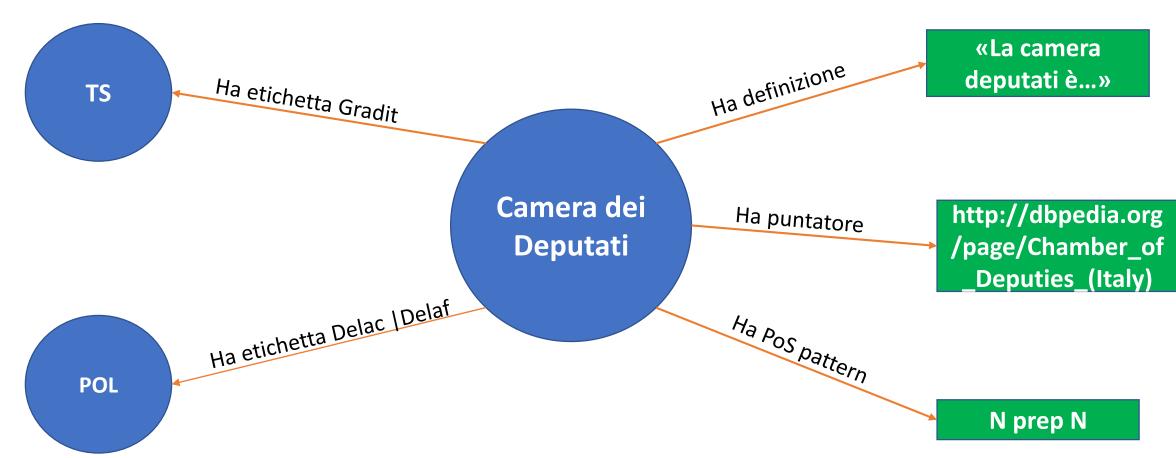
#### Dizionario Ontologico: schema





#### Dizionario Ontologico: schema





PoliSdict (Trotta, Stingo et al., 2018)



Creazione di una risorsa lessicografica digitalizzata per le MWE

Raccolta ed integrazione di diverse fonti informative

Paradigma Open Data

Dizionario elettronico in lingua italiana composto da espressioni polirematiche occorrenti nel parlato spontaneo estratte a partire dal PoliModalCorpus (Trotta, Elia, et al., 2018)

PoliSdict (Trotta, Stingo et al., 2018)



"The extracted MWEs were manually verified using the GRADIT (De Mauro, 2000). This operation has allowed us to identify 356 MWEs compared to 882 identified by DELAC- DELACF and to attribute to each compound expression the respective frequency label documented by the GRADIT"

#### Patrimonio terminologico Network Contacts



Creazione di una risorsa lessicografica digitalizzata per le MWE

Raccolta ed integrazione di diverse fonti informative

Paradigma Open Data

Studio terminologico sul dominio TELCO (telecomunicazioni) per la creazione di risorse linguistiche ausiliarie a sistemi di Information Retrievial e Question Answering

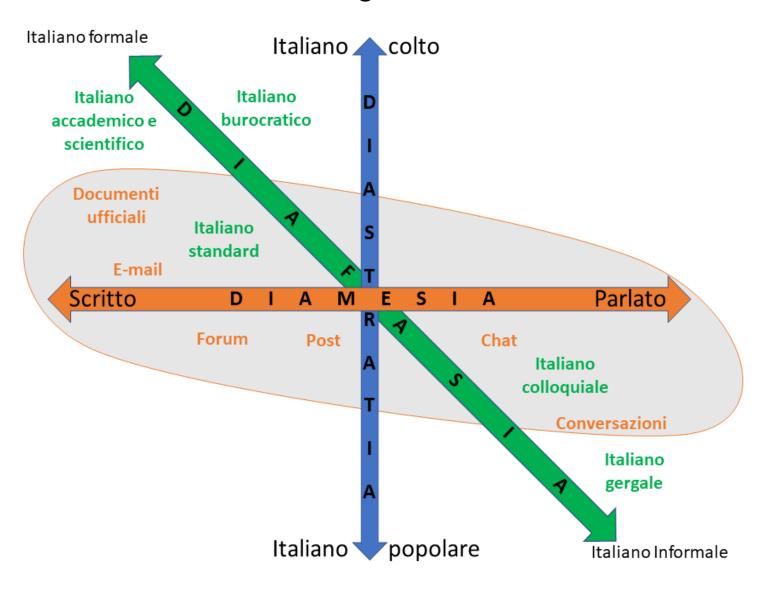
☐ All Inclusive Summer Edition

☐ Scatto alla risposta

☐ Area Clienti

#### Patrimonio terminologico Network Contacts





MWE estratte a partire da un corpus originale, costituito a supporto di soluzioni NLP richiamate da servizi informatici di customer care in area TELCO (telecomunicazioni)

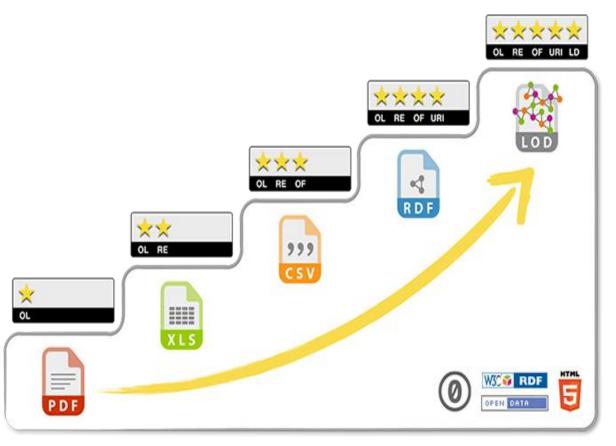
#### Scala d'apertura dei Web Data



Creazione di una risorsa lessicografica digitalizzata per le MWE

Raccolta ed integrazione di diverse fonti informative

Paradigma Open Data



#### Scala d'apertura dei Web Data



Dato reso disponibile in formato machine-readable (.xls).

Dato reso disponibile in formato machine readable non proprietario (.csv, .xml).

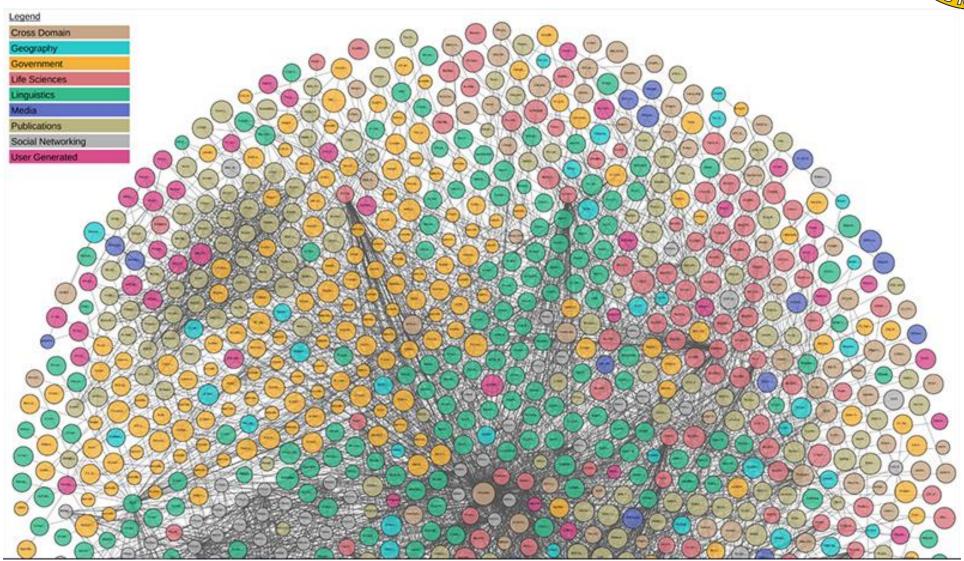
Dato reso disponibile seguendo gli standard e protocolli raccomandati dal W3C. (rdf, owl)

Dato reso disponibile seguendo standard e protocolli W3C e collegato ad altri dataset

Dato reso disponibile sul web in qualsiasi formato (.pdf) e disponibile per il riutilizzo (Open License).

## Linked Open Data Cloud





## Prossimi step



## Opportunità

- Collegamento Dataset di dominio (medico, gastronomico, etc.)
- Potenziamento modello metadescrittivo (acronimi, contesti d'uso, etc.)
- Approccio multilinguistico

#### Criticità

- Implementazione «time consuming»
- Manutenzione periodica



Dizionario Ontologico delle Espressioni Polirematiche dell'Italiano

GRAZIE PER LA VOSTRA ATTENZIONE