

ONTOTEXT
Dal Testo alla Conoscenza per il Web Semantico

Programma di Ricerca FUP-2004
Provincia Autonoma di Trento

**Accesso all'informazione basato su conoscenza:
funzionalità del portale *ONTOTEXT*
(Versione 31/12/2006)**

Luisa Bentivogli, Bernardo Magnini
ITC-irst, Centro per la Ricerca Scientifica e Tecnologica
Via Sommarive 18, Povo (38050) – Trento, Italia
{bentivo, magnini}@itc.it

Technical Report

15 giugno 2006

Indice

Abstract	2
1 Introduzione	3
2 La conoscenza in Ontotext.....	3
3 Il portale dimostrativo ONTOTEXT (versione 31/12/2006).....	4
3.1 Dati disponibili.....	4
3.2 Funzionalità del portale.....	5
4 L'interfaccia utente	6
4.1 Accesso al portale: ricerca di informazioni da parte dell'utente.....	6
4.1.1 Query tramite parole chiave.....	7
4.1.2 Browsing delle categorie ontologiche di Ontotext.....	7
4.2 La risposta del sistema alla query dell'utente	8
4.3 Presentazione dell'informazione.....	11
4.3.1 Informazioni su un'entità.....	11
4.3.2 Informazioni su un argomento (topic)	12
4.3.3 Informazioni su una collocazione	13
4.3.4 Tutti i risultati	14
4.3.5 Accesso ai testi annotati.....	15
5 Prospettive future	16
Bibliografia	16

Abstract

Il progetto Ontotext nasce per studiare e sviluppare tecniche innovative di estrazione della conoscenza al fine di produrre nuova e coerente informazione destinata al Web Semantico. La creazione del portale ONTOTEXT rappresenta lo scenario applicativo del progetto Ontotext, il cui obiettivo finale consiste nel costruire un portale in grado di presentare l'informazione nel modo più possibile flessibile ed aderente alle richieste dell'utente. La principale e innovativa caratteristica del portale ONTOTEXT consiste nel permettere un accesso all'informazione non più basato esclusivamente sui testi, come avviene attualmente nei maggiori portali a disposizione sul Web, bensì basato e guidato dalla conoscenza estratta dai testi stessi.

In questo documento sono presentate le funzionalità del portale dimostrativo di ONTOTEXT che sarà rilasciato al 31/12/2006. Verrà descritta nel dettaglio l'interfaccia utente del portale e analizzati una serie di problemi relativi a come interpretare le domande dell'utente, come trovare le risposte nella base di conoscenza, quali informazioni presentare e come presentarle al meglio.

1 Introduzione

La creazione del portale ONTOTEXT rappresenta lo scenario applicativo del progetto Ontotext, un progetto FU-PAT iniziato nell'ottobre del 2004 che coinvolge ITC-irst e ISTI-CNR. L'obiettivo generale del progetto consiste nello studiare e sviluppare tecniche innovative di estrazione della conoscenza al fine di produrre nuova e coerente informazione destinata al Web Semantico.

L'approccio più comune alla creazione di un Web Semantico si basa sull'annotazione delle risorse disponibili in Rete rispetto a concetti e relazioni definiti all'interno di ontologie, che vengono usate come strumenti concettuali per la comunicazione della conoscenza. L'annotazione semantica delle risorse Web è intesa a facilitarne l'accesso sia da parte di esseri umani sia da parte di agenti artificiali.

L'informazione contenuta nei documenti annotati può essere però incoerente e frammentaria e può cambiare nel tempo. Di conseguenza, il suo sfruttamento può essere complesso, sia da parte di esseri umani che da parte di agenti artificiali. Pur aderendo alla prospettiva del Web Semantico, Ontotext mira a mitigare questo problema integrando nel Web Semantico un nuovo tipo di fonte di informazione risultante da un processo di estrazione di conoscenza guidato da ontologie. Nella nostra prospettiva, la conoscenza contenuta nelle risorse annotate viene estratta e organizzata in una base di conoscenza strutturata, che permette una rappresentazione ed un aggiornamento coerenti dell'informazione, garantendone allo stesso tempo la tracciabilità rispetto alle fonti. I fatti così raccolti possono poi essere impiegati per perfezionare ed estendere le ontologie esistenti. Sia la base di conoscenza che le ontologie estese vengono a loro volta rese disponibili al Web Semantico.

Lo scenario concreto in cui le nuove tecnologie sviluppate verranno testate sarà l'acquisizione automatica di diversi tipi di informazione da articoli del giornale locale *L'Adige*, con il quale è stato concluso un accordo che autorizza l'ITC-irst a processare i dati delle ultime dieci annate del quotidiano da loro fornite. Le potenzialità di integrazione di queste tecnologie saranno dimostrate mediante la realizzazione del portale Web ONTOTEXT, che permetterà ai cittadini di consultare la conoscenza estratta dai testi attraverso un'interfaccia facile da usare e orientata all'utente.

2 La conoscenza in Ontotext

La principale e innovativa caratteristica del portale ONTOTEXT consiste nel permettere un accesso all'informazione non più basato esclusivamente sui testi, come avviene attualmente nei maggiori portali a disposizione sul Web, bensì basato e guidato dalla conoscenza estratta dai testi stessi.

Il progetto Ontotext si pone come obiettivo principale di estrarre e organizzare informazioni su alcune particolari *categorie ontologiche*: entità, oggetti temporali, relazioni, eventi, argomenti (topics) e opinioni.

Le *entità* denotano oggetti o insiemi di oggetti in un dominio di riferimento. Ontotext è interessato alle seguenti entità: "Persone" (es. Paolo Maldini), "Organizzazioni" (es. Vodafone), "Luoghi" (es. Mar Rosso), "Entità Geo-Politiche" (es. Roma). Per ogni entità individuata, viene compilata una *Scheda Ontologica* contenente le principali informazioni estratte per l'entità stessa, ad esempio per le entità di tipo "Persona": Nome, Cognome, Sesso, Professione, Affiliazione ecc.

Una *relazione* è una coppia ordinata di entità, temporalmente collegata (es. il numero di abitanti di una città in un dato anno).

Un *evento* è qualcosa che accade, coinvolge un certo numero di partecipanti e ha come risultato un cambiamento di stato. L'omicidio è un classico esempio di evento.

Gli *argomenti* (topics) sono insiemi di documenti (articoli di giornale nel nostro caso) che riguardano lo stesso avvenimento. Ad esempio "Giro d'Italia 2004", "Dimissioni di Collina". Gli argomenti hanno un certo grado di attivazione in un certo periodo, che dipende da due fattori: (i) la dimensione dell'argomento (il numero dei fatti di cui consiste) e (ii) la frequenza dell'argomento.

Infine, le *opinioni* sono giudizi soggettivi espressi riguardo entità, relazioni, eventi e argomenti. Come primo passo verso l'obiettivo finale di identificare e connotare le opinioni presenti nei testi, l'attività di ricerca in questo ambito è attualmente focalizzata sull'individuazione nei testi dei termini "soggettivi", cioè di quei termini che veicolano non un contenuto oggettivo bensì un'opinione (es. stima vs. cielo) e sul riconoscimento della polarità di questi termini, ovvero se la connotazione del termine è positiva o negativa (es. stima vs. disprezzo).

Per la fine del progetto (30 settembre 2007) verranno studiate le tecnologie e sviluppati gli algoritmi necessari ad estrarre ed organizzare automaticamente la conoscenza relativa a tutte categorie ontologiche di Ontotext descritte sopra. Tale conoscenza sarà estratta dai testi del *corpus Adige-500.000*, un corpus formato dalle ultime dieci annate del quotidiano locale *L'Adige*, per un totale di circa 500.000 articoli. Dati e conoscenza acquisita saranno poi messi a disposizione tramite il portale ONTOTEXT. È prevista dal progetto una fase intermedia di sperimentazione i cui risultati saranno resi disponibili al 31/12/2006 tramite il portale dimostrativo ONTOTEXT che viene descritto nelle sezioni seguenti.

3 Il portale dimostrativo ONTOTEXT (versione 31/12/2006)

Le funzionalità descritte in questo documento si riferiscono al portale intermedio dimostrativo che verrà rilasciato al 31 dicembre 2006 e che darà accesso all'informazione resa disponibile all'interno del progetto per quella data. Nelle due sezioni successive sono presentate le informazioni disponibili e le funzionalità che il portale dimostrativo dovrà avere al 31 dicembre 2006.

3.1 Dati disponibili

Alla data intermedia del 31 dicembre 2006 saranno disponibili in Ontotext i dati descritti di seguito.

Corpora:

- Corpus Adige-500 (I-CAB). 525 articoli (di cui 335 utilizzati come *training set* e 190 come *test set* per gli algoritmi automatici) su cui sono state effettuate una serie di operazioni *manuali*:
 - annotazione di tutte le espressioni temporali
 - annotazione di tutte le entità
 - annotazione di tutte le relazioni di tipo "Personal-Social" sussistenti tra le entità di tipo "Persona", in particolare le relazioni *familiari* (genitore/figlio, ecc.) , *di lavoro* (collegi, capo) e di tipo *personale durevole* (vicini di casa, compagno di classe)
 - per tutte le entità, creazione di una *Scheda Ontologica* contenente una serie di informazioni rilevanti per quella data entità
- Corpus Adige-50.000. 52.000 articoli annotati automaticamente con informazione di tipo linguistico (lemma e parte del discorso), entità nominate ed espressioni

temporali. Su questo corpus di medie dimensioni vengono attualmente testate tutte le procedure automatiche prima di passare al corpus finale.

- Corpus Adige-500.000. È il corpus finale di Ontotext, composto di 500.000 articoli. Al 31/12/2006 sarà annotato con informazione di tipo linguistico (lemma e parte del discorso), entità nominate ed espressioni temporali mentre per la fine del progetto tutte le procedure funzioneranno su questo corpus che sarà interamente accessibile tramite il portale.

Ontologie:

- Ontologia delle Persone
- Ontologia delle Organizzazioni
- Ontologia dei Luoghi
- Ontologia delle Entità Geo-Politiche

Procedure automatiche:

- Processori linguistici per l'annotazione automatica dei diversi tipi di informazione linguistica necessaria in Ontotext (lemma, parte del discorso, collocazioni, sensi di parola)
- Sistema di Named Entity Recognition per il riconoscimento delle entità nominate nei testi
- Sistema di riconoscimento e normalizzazione delle espressioni temporali nei testi
- Algoritmi per il riconoscimento di argomenti (topics) e classificazione dei testi in base all'argomento di cui trattano
- Algoritmi per il riconoscimento di opinioni positive e negative nei testi

Come vedremo nel dettaglio in seguito, alcune funzionalità del portale dimostrativo accederanno al corpus I-CAB e quindi alle informazioni codificate manualmente nei testi o estratte manualmente dai testi, mentre altre sfrutteranno già il corpus finale Adige-500.000 e le informazioni estratte tramite le procedure automatiche dai 500.000 testi.

Dal punto di vista delle categorie ontologiche studiate in Ontotext, al 31 dicembre 2006 non saranno a disposizione informazioni sugli eventi.

3.2 Funzionalità del portale

Al 31/12/2006 saranno disponibili dal portale dimostrativo le seguenti funzionalità:

- 1) Citografo
- 2) Ricerca di una o più parole
- 3) Ricerche avanzate (parole in una determinata parte del discorso, lemmi, lemmi in una determinata parte del discorso, sottostringhe, ricerca di più parole con restrizioni sulla modalità di ricerca)
- 4) Ricerca di collocazioni
- 5) Accesso alle schede informative delle collocazioni
- 6) Ricerca/browsing di entità
 - a) entità nominate
 - b) menzioni di entità
 - c) co-riferenze di entità a livello di singolo documento
 - d) co-riferenze di entità a livello di corpus
- 7) Accesso alle schede ontologiche delle entità
- 8) Ricerca/browsing di argomenti

- 9) Ricerca/visualizzazione di relazioni tra persone
- 10) Ricerca di opinioni
- 11) Accesso ai testi e ai diversi livelli di annotazione

Le funzionalità di ricerca/browsing di menzioni di entità (5b) e di co-referenza (5c e 5d), di accesso alle schede ontologiche delle entità (6) e alle relazioni tra persone (8) sono disponibili solo per i dati ricavati *manualmente* dal corpus I-CAB, e quindi su un totale di 500 documenti. Tutte le altre funzionalità utilizzano invece informazioni estratte *automaticamente* dal corpus Adige-500.000.

4 L'interfaccia utente

Data la conoscenza a disposizione, unitamente ai testi da cui è stata estratta, la sfida di Ontotext consiste nel costruire un portale in grado di presentare tale conoscenza nel modo più possibile flessibile ed aderente alle richieste dell'utente. In questa sezione descriveremo in dettaglio l'interfaccia utente del portale ONTOTEXT e analizzeremo una serie di problemi relativi a come interpretare le domande dell'utente, come trovare le risposte nella base di conoscenza, quali informazioni presentare e come presentarle al meglio.

4.1 Accesso al portale: ricerca di informazioni da parte dell'utente

L'utente ha la possibilità di effettuare due tipi di ricerca:

- (a) eseguire una ricerca standard, tipo Google, inserendo una o più parole nell'apposito campo
- (b) navigare direttamente le categorie ontologiche di Ontotext

La pagina di accesso al portale, da cui sono disponibili contemporaneamente le due modalità di ricerca, è riportata in Figura 1.

PORTALE WEB ONTOTEXT

Cerca nel giornale "L'Adige":

Inserisci una o più parole

[Ricerca avanzata](#)

Cerca una categoria ontologica di Ontotext nel giornale "L'Adige":

<input type="button" value="PERSONE"/>	<input type="button" value="ORGANIZZAZIONI"/>
<input type="button" value="LUOGHI"/>	<input type="button" value="ENTITÀ GEO-POLITICHE"/>
<input type="button" value="EVENTI"/>	<input type="button" value="ARGOMENTI"/>

Figura 1: pagina di accesso al portale Ontotext

4.1.1 Query tramite parole chiave

Nella prima modalità di ricerca, l'utente inserisce una o più parole nell'apposito campo. L'utente può fare vari tipi di query, ad esempio una domanda vera e propria: "qual è attualmente il sindaco di Trento?"; oppure inserire una serie di parole chiave: "Elenco iscritti partito verdi"; oppure essere ancora più generico e inserire semplicemente una parola: "margherita". È il sistema che si occupa poi di interpretare la query dell'utente in modo da restituire le informazioni più aderenti alla richiesta.

È inoltre prevista una funzionalità di ricerca avanzata in cui l'utente può ad esempio cercare sottostringhe di parole o fare ricerche più raffinate dal punto di vista linguistico, come cercare una parola in una certa parte del discorso o un lemma, in una determinata parte del discorso o meno. Questa modalità è rappresentata in Figura 2.

Ricerca Avanzata

Cerca la frase intera		<input type="text"/>
Cerca tutte le parole		<input type="text"/>
Cerca una qualunque delle parole		<input type="text"/>
Cerca una sottostringa		<input type="text"/>
cerca una parola*	<input type="text" value="nome"/>	<input type="text"/>
Cerca un lemma **	<input type="text" value="nome"/>	<input type="text"/>

***Parola in diverse parti del discorso:**
dato: participio passato verbo "dare"
dato: aggettivo
dato: nome

**** Lemma:** voce registrata in un dizionario; ne introduce la definizione
Lemma in diverse parti del discorso ardire: verbo
ardire: nome

Figura 2: modalità di ricerca avanzata

4.1.2 Browsing delle categorie ontologiche di Ontotext

Nella seconda modalità di ricerca, l'utente può fare direttamente il browsing delle categorie di Ontotext. Ogni categoria è strutturata in una gerarchia e l'utente può navigarla, scegliendo le sottocategorie a cui è interessato. In alternativa, può ignorare la gerarchia e visualizzare l'elenco completo in ordine alfabetico degli appartenenti alla categoria selezionata. Una volta selezionata una determinata entità o argomento, si raggiunge direttamente la relativa scheda informativa descritta in dettaglio nelle sezioni 4.3.1 e 4.3.2.

Ad esempio, come mostrato in Figura 3, l'utente interessato alle "organizzazioni" di tipo "istruzione" può cliccare sulla sottocategoria nella parte destra della pagina e nella parte sinistra compaiono tutte le entità di quel tipo. Cliccando su una di queste si raggiunge la scheda che raccoglie tutte le informazioni a disposizione del sistema su quella entità.

Lista delle Organizzazioni	
Scegli una lettera	Scegli un tipo
Tipo: istruzione	<ul style="list-style-type: none"> • <u>tutte le organizzazioni</u> <ul style="list-style-type: none"> • <u>commerciale</u> • <u>istruzione</u> • <u>governativa</u> • <u>scienza medica</u> • <u>sportiva</u> • ...
<ul style="list-style-type: none"> • Dipartimento di Ingegneria civile e ambientale - Università di Trento • Centro Universitario per la Difesa Idrogeologica dell' Ambiente Montano (CUDAM) - Università di Trento • <u>Università di Trento</u> • Museo Civico di Storia Naturale - Verona • ATNA - Trento • Conservatorio "Monteverdi" - Bolzano • Fondazione Concorso Pianistico Internazionale "Busoni" • Accademia di Danza Orientale Margarita - Trento • scuola professionale di sartoria - Rovereto • Istituto d' arte del tessuto - Rovereto • scuole medie di Riva del Garda • Istituto Luce • Biennale - Venezia 	

Figura 3: funzionalità di browsing delle categorie di Ontotext

4.2 La risposta del sistema alla query dell'utente

Nel caso in cui l'utente abbia deciso di fare una query cercando parole chiave, il sistema interpreta la query e fornisce una risposta. I risultati ottenuti vengono raggruppati dal sistema in maniera intelligente:

- (i) viene visualizzato l'elenco delle *entità* che soddisfano la query.
- (ii) viene visualizzato l'elenco degli *argomenti* che soddisfano la query.
- (iii) viene presentato l'elenco delle *collocazioni*¹ contenenti la/e parola/e della query. Le collocazioni rappresentano un livello intermedio tra la conoscenza e il testo puro.
- (iv) vengono presentati tutti i risultati ottenuti, senza tener conto della conoscenza del sistema. La query viene trattata come un insieme di parole che vengono cercate nei documenti. Il risultato viene visualizzato o sotto forma di *concordanze*² (quando viene trovata esattamente la parola o la frase inserita dall'utente) o di *snippets*³ di documenti (quando nei documenti vengono trovate separatamente le singole parole che compongono la query e non la query stessa così come è stata digitata).

¹ Una collocazione è una sequenza di parole che co-occorrono abitualmente e il cui significato complessivo non è completamente ricostruibile dal significato delle singole parole che la compongono. Ad esempio "senso unico", "montagne russe".

² Per concordanza si intende la lista di tutte le occorrenze di una data parola in un testo, assieme al contesto in cui quella parola occorre. La parola cercata è centrata rispetto alla riga ed evidenziata graficamente, il contesto è rappresentato da un certo numero di parole che precedono e che seguono la parola cercata.

³ Uno snippet è l'insieme delle porzioni del testo in cui compaiono le parole della query, riunite in un unico paragrafo ai fini della presentazione all'utente.

Tutte queste informazioni vengono recuperate dal corpus Adige-500 e sono il risultato di annotazioni automatiche dei testi.

Vediamo ad esempio in Figura 4 il risultato fornito dal sistema alla query costituita da una singola parola: “margherita”. Due *entità* di tipo diverso (Persone e Organizzazioni) e un *argomento* soddisfano la query; oppure l’utente ha la possibilità di vedere tutti i risultati assieme. In questo esempio il sistema non ha trovato nel corpus nessuna collocazione in cui figurasse la parola “margherita”, quindi il livello delle collocazioni non viene presentato.

The image shows a search interface. At the top, there is a search bar with the text "Inserisci una o più parole:" and a search button labeled "Cerca". The search bar contains the word "margherita". To the right of the search bar is a link labeled "Ricerca avanzata". Below the search bar, there is a green box containing the following text:

Risultati per *margherita*:

- margherita - Persona
 - Margherita Cogo
- margherita - Organizzazione
 - La Margherita
- margherita – Argomento
 - festa della Margherita - Polignano

Vedi tutti i risultati per **MARGHERITA**

Figura 4: risposta del sistema alla query dell’utente “margherita”

Nell’esempio in Figura 5, data la query “rosso” il sistema presenta sia entità (di tipo Persone e Luoghi) sia una serie di collocazioni di cui la parola “rosso” fa parte.

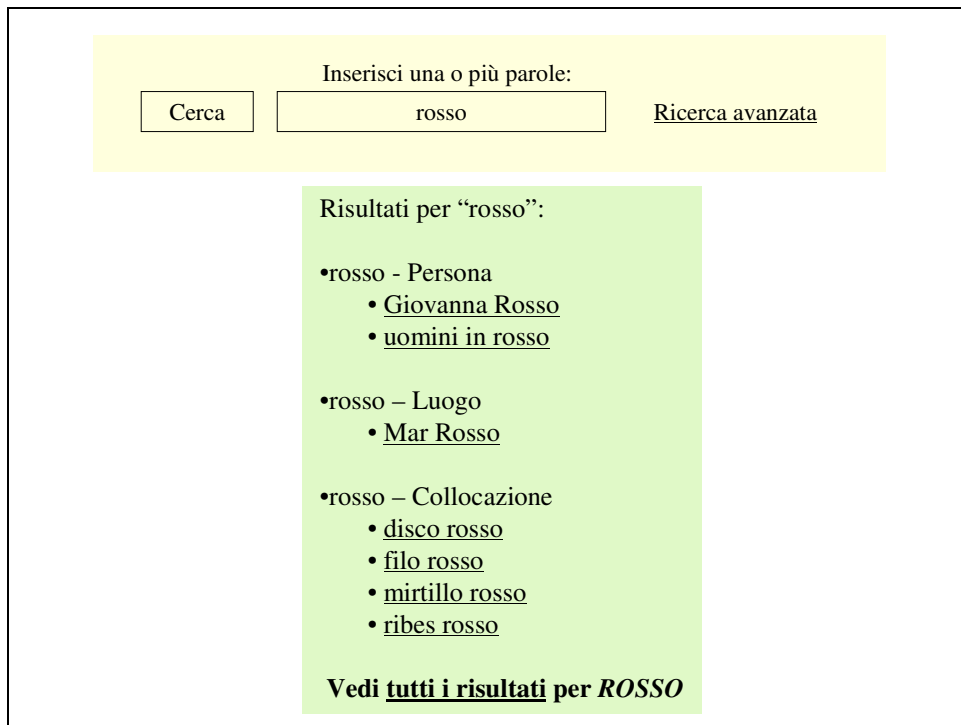


Figura 5: risposta del sistema alla query dell'utente "rosso"

Infine, nell'esempio in Figura 6 il sistema interpreta la query complessa dell'utente e riesce a trovare un'entità di tipo Organizzazioni (il partito politico dei verdi) che la soddisfa. Restituisce inoltre un elenco di documenti in cui compaiono le singole parole che compongono la query.

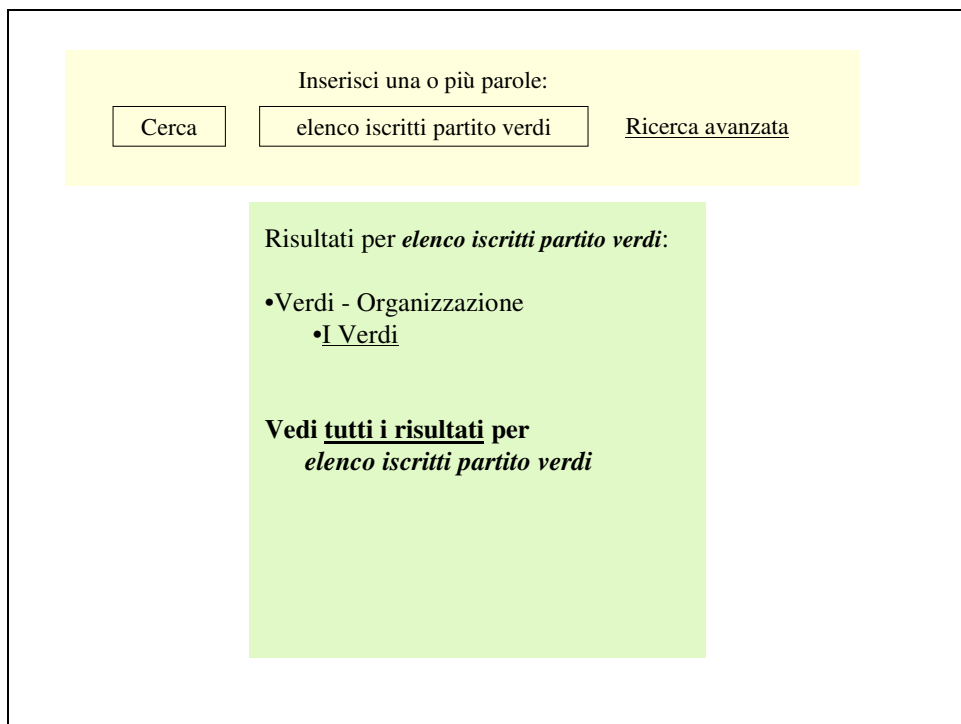


Figura 6: risposta del sistema alla query dell'utente "elenco iscritti partito verdi"

4.3 Presentazione dell'informazione

Sia nel caso di una query a parole chiave (cfr. sezioni 4.1.1 e 4.2) sia nel caso di browsing diretto delle categorie ontologiche di Ontotext (cfr. sezione 4.1.2.), il sistema presenta all'utente tutta la conoscenza a disposizione su entità ed argomenti, fornendo inoltre l'accesso ai singoli documenti annotati. Inoltre, in caso di query a parole chiave, sono attive anche la funzione di visualizzazione di collocazioni e quella di "concordancer", che dà all'utente la possibilità di visualizzare tutti i documenti che contengono le parole della sua query, indipendentemente dalle conoscenze possedute dal sistema. Nelle sezioni che seguono, le informazioni presentate e le modalità di presentazione vengono descritte dettagliatamente.

4.3.1 Informazioni su un'entità

L'utente clicca su una determinata entità e il sistema restituisce una scheda che contiene tutta la conoscenza posseduta su quella entità. La schermata è suddivisa in quattro frame contenenti:

1. scheda ontologica dell'entità. La scheda ontologica di un'entità contiene le principali informazioni estratte per l'entità stessa, contestualizzate temporalmente. Ad esempio per le persone sono presenti informazioni come Nome, Cognome, Sesso, Professione, Affiliazione ecc. e, dove necessario come ad esempio per Professione e Affiliazione, è presente anche l'informazione del periodo a cui quell'informazione si riferisce. Ogni valore degli attributi della scheda ontologica è cliccabile e rimanda all'articolo da cui quel valore è stato estratto.
2. le concordanze dell'entità, ordinate dall'articolo più recente a quello meno recente. Ogni concordanza è collegata all'articolo da cui è stata estratta. In questo frame è disponibile anche la funzionalità di visualizzazione delle opinioni positive e negative sull'entità in esame, contestualizzate temporalmente. Se l'entità è una persona, è possibile avere anche le opinioni (positive e negative) espresse da quell'entità, sempre temporalmente contestualizzate.
3. accesso al citografo relativo all'entità. Il citografo è un grafico che mostra l'andamento temporale delle citazioni dell'entità presa in considerazione nei vari articoli di giornale.
4. possibilità di visualizzare altre entità collegate tramite una relazione all'entità in esame e di navigare la porzione di ontologia in cui si trova l'entità stessa.

I frame 1 e 4 contengono informazione estratta manualmente dal corpus I-CAB e sono quindi disponibili solo ed esclusivamente se l'entità presa in considerazione è presente nel corpus I-CAB.

Per quanto riguarda il frame 2, se nel corpus I-CAB è presente l'entità presa in considerazione l'utente può scegliere se visualizzare le *concordanze dell'entità* nel corpus I-CAB dove le co-riferenze a livello di corpus sono risolte, oppure se visualizzare le *concordanze a livello di Entità Nominate* estratte automaticamente dal corpus Adige-500.000, dove la co-riferenza delle entità a livello di corpus non è stata codificata. Se l'entità non è presente in I-CAB saranno disponibili solo le concordanze di Entità Nominate estratte da Adige-500.000.

Anche il frame 3 è sempre disponibile ma, poiché il citografo necessita di una grossa quantità di dati e quindi lavora su Adige-500.000 dove le co-riferenze delle entità non sono state risolte, non funziona a livello di entità bensì di occorrenza di Entità Nominata.

Nell'esempio in Figura 7 è mostrata la scheda personale che il sistema presenta all'utente che ha cliccato sull'entità di tipo persona "Margherita Cogo" (cfr. Figura 4).

Scheda Personale di Margherita Cogo

Scheda Ontologica		Occorrenze di Margherita Cogo nel giornale	
Nome	Margherita	<ol style="list-style-type: none"> 1. Così il congresso condiziona il lavoro dell'esecutivo. Si accendono anche i Ds. Cogo critica, poi l'esecutivo frena le polemiche. Sei - righe - sei per dire che questa volta va bene, che si può andare avanti. 2. Peccato che questo succinto comunicato segua di poche ore la pubblicazione di una intervista ad una Margherita Cogo mai così agguerrita. 3. Macché, ora è Margherita Cogo. 4. E non era quest'ultima a guidare il partito dei filo - governativi, dei fedelissimi del presidente disposti a chinare il capo pur di non mollare le loro poltrone vellutate? 5. se la Cogo prende posizione contro Dellai, è evidente la sua preoccupazione di occupare lo spazio di chi crede che la sinistra sia troppo appiattita sulla Margherita, spazio ora occupato da Andreoli e Chiodi. 	<ol style="list-style-type: none"> 1. « Quando io dicevo le cose che ora sostiene la Cogo, passavo per quello che vuole spaccare tutto.
Cognome	Cogo		
Soprannome			
Titolo			
Sesso	Femmina		
Attività		<ol style="list-style-type: none"> 2. Cultura. Venti milioni di euro, « spalmati » solo su sei delle 87 realtà museali pubbliche e private del Trentino Musei, spese altissime. Vertice della Cogo con Tirolo e Alto Adige in confronto Trento ha conti stratosferici. Di CORONA PERER Venezia spende circa 10 milioni di euro per i musei di tutta la città, la Provincia di Bolzano " spalma " i suoi 8 milioni e rotti su 8 musei provinciali e altre 70 realtà private, il Tirolo non ne spende più di 9, « (nota bene) il suo budget negli ultimi 20 anni ha avuto incremento del 400%. 3. « Sì, la cifra è un po' imbarazzante » ha esordito la responsabile delle politiche culturali Margherita Cogo all'incontro convocato al Mart di Rovereto con le colleghe Sabina Kastlatter Mur, assessore alla cultura della Provincia di Bolzano e l'omologa Elisabeth Zenon per il Land Tirolo. 4. Dopo aver sottolineato che il dato può essere letto come la prova dell'importanza che il Trentino dà alla cultura, l'assessore Cogo ha ammesso che la cifra può essere anche indicativa di una " gestione dispendiosa ". 5. Quanto al personale, Margherita Cogo ha fatto capire che il problema, certamente non secondario, è per certi versi elico. 	
Affiliazione	Comune di Tione (1985-1998) Regione TAA (1999-2002) Provincia TN (2003 - oggi)		
Ruolo	•Consigliere (1985 - 1993) •Sindaco (1993-1998) •Presidente (1999-2002) •Assessore alla cultura (2003-oggi) •Vicepresidente (2003 - oggi)		
Provenienza	Tione (Trento)		
Miscellanea	DS		
		VEDI OPINIONI "DI" E "SU" MARGHERITA COGO	

Citazioni di Margherita Cogo nel giornale	Altre entità collegate:	Ontologia di Margherita Cogo
	<ul style="list-style-type: none"> -Persone <li style="padding-left: 20px;">-Lorenzo Dellai -Organizzazioni <li style="padding-left: 20px;">-Provincia di Trento -DS 	

Figura 7: Scheda informativa di un'entità

4.3.2 Informazioni su un argomento (topic)

L'utente clicca su un determinato argomento e il sistema restituisce una scheda strutturata in modo parallelo a quella per le entità. Anche la scheda per gli argomenti è suddivisa in quattro campi contenenti:

1. la scheda informativa dell'argomento contenente le informazioni principali sull'argomento stesso
2. l'elenco degli articoli che parlano dell'argomento
3. il grafico che riporta l'andamento temporale dell'attivazione dell'argomento
4. la lista di altri argomenti correlati all'argomento in esame e la relativa ontologia

La scheda informativa sugli argomenti utilizza informazioni estratte automaticamente da Adige-500.000.

Nell'esempio in Figura 8 è mostrata la scheda che il sistema presenta all'utente che abbia cercato informazioni sull'argomento "strage di Beslan".

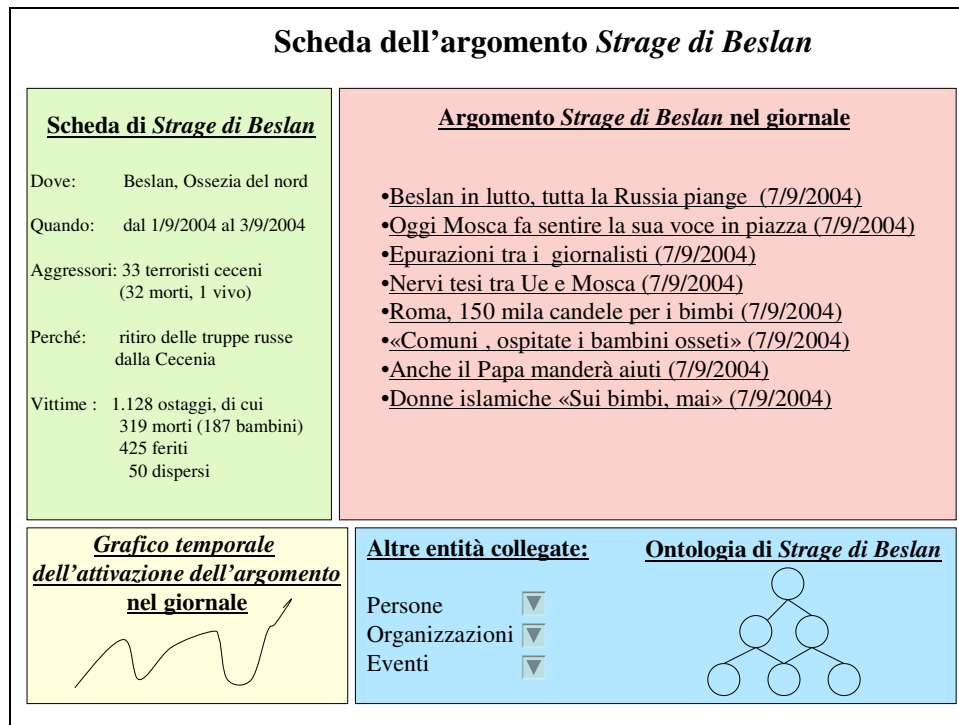


Figura 8: Scheda informativa di un argomento

4.3.3 *Informazioni su una collocazione*

Quella sulle collocazioni è un'informazione di tipo linguistico che viene presentata all'utente in una apposita scheda strutturata in modo parallelo a quella per le entità. Anche la scheda per le collocazioni è suddivisa in quattro campi contenenti:

1. la definizione della collocazione (tratta dal database lessicale MultiWordNet o da un dizionario)
2. le concordanze della collocazione. Ogni concordanza è collegata all'articolo da cui è stata estratta.
3. accesso al citografo relativo alla collocazione
4. le relazioni semantiche relative al concetto e la porzione di gerarchia (navigabile) di MultiWordNet in cui si trova il concetto.

La scheda informativa sulle collocazioni utilizza informazioni estratte automaticamente da Adige-500.000.

Nell'esempio in Figura 9 è mostrata la scheda che il sistema presenta all'utente che abbia cliccato sulla collocazione "ribes rosso" (cfr. Figura 5).

Scheda di *ribes rosso*

Scheda lessicografica

Definizione di *ribes rosso*
in MultiWordNet

The word "ribes_rosso" has 1 senses:

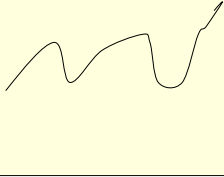
Noun	
1. Ribes_rubrum, ribes, ribes_rosso	(Botany) [cultivated European current bearing small edible red berries]

Occorrenze di *ribes rosso* nel giornale

Occurrences of **ribes rosso**: 1 (1 sentences / 1 files)

1. ed il nero, il **ribes rosso** e nero, la ciliegia

Citazioni di *Ribes rosso* nel giornale

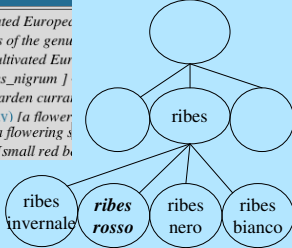


Altri concetti collegati:

Noun

- 1. Ribes_rubrum, ribes, ribes_rosso -- (Botany) [cultivated European
- > 1. ribes -- (Botany) [any of various deciduous shrubs of the genu
- => 1. Ribes_rubrum, ribes, ribes_rosso -- (Botany) [cultivated Eur
- => 1. [black_currant, European_black_currant, Ribes_nigrum]
- => 1. [white_currant, Ribes_sativum] -- (Botany) [garden curra
- => 1. [winter_currant, Ribes_sanguineum 1 -- (Botany) [a flower
- IS MEMBER-OF: 1. [Ribes, genus_Ribes] -- (Biology) [a flowering s
- HAS PART: 1. [red_currant] -- (Botany, Gastronomy) [small red b

Ontologia di *Ribes rosso*



```

graph TD
    A((ribes)) --- B((ribes invernale))
    A --- C((ribes rosso))
    A --- D((ribes nero))
    A --- E((ribes bianco))
  
```

Figura 9: scheda informativa di una collocazione

4.3.4 Tutti i risultati

Come già visto nella sezione 4.2, in caso di query a parole chiave, è attiva la funzionalità “tutti i risultati”, che dà all’utente la possibilità di visualizzare tutti i documenti che contengono le parole della sua query, indipendentemente dalle conoscenze possedute dal sistema.

Sono previste due modalità di visualizzazione di tutti i risultati della query. Se il sistema trova esattamente la parola o la frase inserita dall’utente, vengono restituiti i documenti sotto forma di *concordanze* di quella parola o frase. Se invece nei documenti vengono trovate le singole parole che compongono la query e non la query stessa così come è stata digitata il sistema restituisce *snippets* di documenti. La prima modalità è esemplificata in Figura 10, mentre la seconda in figura 11.

Tutti i risultati per la parola **ROSSO**

Occurrences of **rosso**: 12 (11 sentences / 8 files)

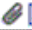


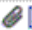







1.	  	avvarranno delle voci di Giovanna Rosso e Claudia Tomasini della filodrammatica
2.	  	è un cartello bianco e rosso con la scritta : qui
3.	  	Il tutto accompagnato dal filo rosso della musica che attraversa i
4.	  	in perdita Trentofrutta Il < rosso > di 355.000 euro dovuto
5.	  	il nero , il ribes rosso e nero , la ciliegia la mora , il mirtillo rosso ed il nero , il
6.	  	dato parere positivo ; disco rosso invece per la zona di
7.	  	località balneari egiziane sul Mar Rosso , vicine a Sharm el  stazione balneare egiziana sul Mar Rosso , al confine con lo  Notte di terrore sul Mar Rosso . Tre esplosioni , in  nel Sinai e sul Mar Rosso > .
8.	  	in favore degli uomini in rosso apparirebbe pure giustificato , ecco

Figura 10: risultato della query per “rosso” (concordanza)

All results for word “elenco iscritti partito verdi”

Elenco deputati eletti · **Elenco** senatori eletti · Elezioni politiche 2006 ... Gli ultimi candidati **iscritti** ... **Partito: VERDI**. Scrivi una mail al candidato ...

Hanno poi preso la parola alcuni **iscritti** del **partito**. ... il **partito** ha rivolto poca attenzione agli inviti di altre forze politiche dell'Unione (**Verdi**, ...

L'**elenco** dei senatori divisi per gruppi. (**Elenco** parlamentari 4.5.2006) ... di Gruppi con meno di dieci **iscritti**, purchè rappresentino un **partito** o un ...

Figura 11: risultato della query per “elenco iscritti partito verdi” (snippets)

4.3.5 Accesso ai testi annotati

Una caratteristica fondamentale in Ontotext è la tracciabilità della conoscenza estratta automaticamente dai testi. Per questo motivo da ogni informazione fornita è possibile risalire alle fonti da cui è stata estratta. All'utente viene presentato l'articolo di giornale con evidenziata la porzione di informazione tramite cui si è arrivati a quel testo. È inoltre possibile visualizzare i diversi tipi di annotazione presenti nel documento: annotazione linguistica (sempre automatica) e annotazione semantica (manuale per il corpus I-CAB e automatica per il corpus finale Adige-500.000).

Da questa pagina è inoltre disponibile la funzionalità “mostra i testi simili” tramite la quale si accede direttamente ad altri articoli che parlano dello stesso argomento.

La Figura 12 mostra un esempio di testo annotato raggiunto partendo dalla scheda informativa di Margherita Cogo (cfr. Figura 7) in cui sono evidenziate tutte le occorrenze dell'entità di tipo Persona “Margherita Cogo”.

Linguistic annotation: Noun Verb Adjective Adverb Multiword

Semantic annotation: Time Mention una Margherita Cogo mai così a... Select ORGANIZATIONS » show

l'Adige - training/adige20040907_id405610

A sinistra La Quercia soddisfatta . Così il congresso condiziona il lavoro dell' esecutivo

Si accodano anche i

Ds [Cogo] critica , poi l' esecutivo frena le polemiche Sei - righe - sei per

dire che questa volta va bene , che si può andare avanti . Ecco il comunicato diffuso ieri sera dopo una breve riunione : « L' esecutivo dei Ds del Trentino prende atto positivamente del nuovo metodo di confronto politico e programmatico con cui la coalizione di centrosinistra autonomista ha ripreso l' attività dopo la pausa estiva . In particolare , l' esecutivo dei Ds prende atto positivamente del fatto che la coalizione si riunirà nei prossimi giorni per discutere e approfondire i temi legati alla ricerca , alla mobilità , alla riforma istituzionale e all' energia » . Insomma , avanti così . Peccato che questo succinto comunicato segua di poche ore la pubblicazione di una intervista ad [una Margherita Cogo mai così agguerrita] . Sì , con tanto di dichiarazioni bellicose per il rispetto dei programmi e in difesa dell' « idea » che ha la sinistra della riforma istituzionale . I Democratici di sinistra del Trentino sono sempre in grado di stupire :

non appena ti crei una « fotografia » della Quercia , ecco che questa si dissolve .

Pensi di aver piazzato tutti i nomi nelle diverse categorie e - in un battibaleno - il castello salta per aria .

Davvero un rebus senza fine .

Ma non era Mauro Bondi quel « Follini de noantri » che richiamava al rispetto dei programmi e avanzava critiche sulla riforma istituzionale e sul disegno di riordino della ricerca ?

Macché , ora è [Margherita Cogo] .

E non era [quest' ultima] a guidare il partito dei filo - governativi , dei fedelissimi del presidente disposti a chinare il capo pur di non mollare le loro poltrone vellutate ?

No , niente vero .

Ora si resta in attesa di nuove categorie , nuove prese di posizione , ulteriori cordate .

Certo che il congresso della Quercia condiziona un bel po' la politica trentina , anche se ha ragione l' onorevole Gigi Olivieri quando dice che tutto sommato « siamo un circolo di mille iscritti , e per questo dobbiamo guardare oltre il partito » .

Sì , perché molti degli interventi di questi giorni possono essere letti come « aggiustamenti » in vista del rinnovo della segreteria : se [la Cogo] prende posizione contro Dellai , è evidente la [sua] preoccupazione di occupare lo spazio di chi crede che la sinistra sia troppo appiattita sulla Margherita , spazio ora occupato da Andreolli e Chiodi .

Il più ironico in questi giorni pare essere Mauro Bondi , che si sta beccando qualche riconoscimento :

« Quando io dicevo le cose che ora sostiene [la Cogo] , passavo per quello che vuole spaccare tutto .

Figura 12: visualizzazione di un testo annotato

5 Prospettive future

Il progetto Ontotext è in corso di sviluppo. Per la fine del progetto sono previste ulteriori attività. Il corpus Adige-500.000 sarà interamente processato per estrarre tutta la conoscenza prevista in Ontotext, eventi compresi. Il database dei fatti (attualmente composto dalle informazioni estratte manualmente dal corpus I-CAB) sarà enormemente arricchito con la conoscenza estratta automaticamente dai 500.000 testi del corpus Adige-500.000. Il portale Ontotext rilasciato in versione dimostrativa per il 31 dicembre 2006 sarà aggiornato e arricchito durante tutto lo svolgimento del progetto fino ad arrivare alla versione definitiva con la fine del progetto stesso.

Bibliografia

L. Bentivogli, E. Pianta. Detecting hidden multiwords in bilingual dictionaries. In *Proceedings of the tenth EURALEX International Congress*, Copenhagen, Denmark, August 14-17, 2002, pp. 785-793

A. Esuli and F. Sebastiani. 2006. Determining Term Subjectivity and Term Orientation for Opinion Mining. In *Proceedings of the 11th Conference of the European Chapter of the Association for Computational Linguistics (EACL 2006)*.

B. Magnini, E. Pianta, O. Popescu, M. Speranza: Ontology Population from Textual Mentions: Task Definition and Benchmark, to appear in *Proceedings of the OLP-2 workshop, 2nd workshop on Ontology Learning and Population*, Sydney, Australia, 2006.

B. Magnini, E. Pianta, C. Girardi, M. Negri, L. Romano, M. Speranza, V. Bartalesi-Lenzi, and Rachele S.i. I-CAB: the Italian Content Annotation Bank. In *Proceedings of LREC 2006 - 5th Conference on Language Resources and Evaluation, 22-28/5/2006, Genova (Italy)*.

Magnini B., Negri M., Pianta E., Romano L., Speranza M., and Sprugnoli R. From Text to Knowledge for the Semantic Web: the ONTOTEXT Project. In *Proceedings of SWAP 2005, Semantic Web Applications and Perspectives*, Trento, Italy, 14-15-16 December, 2005.

E. Pianta, L. Bentivogli, C. Girardi, B. Magnini. Representing and Accessing Multilevel Linguistic Annotation using the MEANING Format. In *Proceedings of the EACL-06 Workshop on Multi-dimensional Markup in Natural Language Processing (NLPXML-2006)*, Trento, Italy, April 4, 2006.

Saquete E., Martinez-Barco P., Munoz R., Negri M., Speranza M., and Sprugnoli R. Multilingual Extension of a Temporal Expression Normalizer using Annotated Corpora. In *Proceedings of the EACL 2006 Workshop on Cross-Language Knowledge Induction*, Trento, Italy, April 3, 2006.

Siti Web

- <http://www.elpais.es/> (quotidiano El Pais)
- <http://www.google.com> (motore di ricerca GOOGLE)
- <http://vivisimo.com> (Vivisimo, clustering automatico di pagine web)