



Quality-driven Evolution in Information Integration Systems

Ph.D. in Computer Science Course - XXVIII Series

Riccardo Porrini

Supervisor: Dott. Palmonari

Tutor: Prof. Messina

shoppydoo

La tua guida allo shopping online

Like 21k

Segui

Accedi | Registrati

ShoppyDoo > Telefonia > Cellulari > Offerte iPhone 6

Risultati per: **iphone 6** in **Cellulari**

Aggiungi la ricerca ad una lista

Cerca **IPHONE 6** anche in:

Accessori Cellulari Cover per Cellulari Accessori Fitness Ricambi per Cellulari Altre categorie

Filtra la ricerca

PREZZO

€ 192,02 - € 1347,06

MARCA

Apple (534)

SISTEMA OPERATIVO

iOS (430)

TIPO DI DISPOSITIVO

Smartphone (504)

Phablet (157)

PROCESSORE (CPU)

Dual-core (2)

CONNESSIONI

LTE (4G) (427)

NFC (491)

3G (3)

Bluetooth (75)

Wi-Fi (74)

ALTRE FUNZIONI

Con GPS Integrato (504)

Touch screen (504)

Risultati: **545**

Ordina per: Prezzo Nome Popolarità



Apple iPhone 6 16GB

[Recensione](#) | [Scheda tecnica](#)

Versione del nuovo iPhone 6 con schermo Retina HD da 4,7 pollici e 16 GB di memoria integrata. E' realizzato in alluminio anodizzato, acciaio e vetro, con uno spessore di soli 6,9 mm. Al suo interno, si trova un potente ed efficiente processore A8. Supporta la connessione LTE ed la fotocamera iSight da 8 Mpx si avvale della tecnologia Focus Pixel, che rende l'autofocus ancor più veloce e

da € 569

Vai al più economico

Confronta i prezzi

in **41** negozi

Confronta prodotto

Aggiungi a una lista

Prezzo desiderato



Apple iPhone 6 64GB

[Recensione](#) | [Scheda tecnica](#)

Del nuovo iPhone con display Retina HD da 4,7 pollici, questa è la versione con 64 GB di memoria integrata. Supporta, come gli altri modelli, la veloce connessione LTE ed è dotato del nuovo processore A8. Più potente ed efficiente del precedente, migliora i tempi di autonomia. La fotocamera iSight da 8 Mpx è ora dotata di tecnologia Pixel Focus, veloce e precisa. Nuovo e,

da € 640

Vai al più economico

Confronta i prezzi

in **35** negozi

Confronta prodotto

Aggiungi a una lista

Prezzo desiderato

categories

facets

Country of Origin

USA (320)

France (91)

Italy (40)

Spain (18)

Australia (17)

Bulgaria (10)

Chile (8)

+ See more

Vintage

☐ No Vintage (96)

☐ 2013 (72)

☐ 2012 (102)

☐ 2011 (84)

☐ 2010 (50)

+ See more

1-24 of 8,933 results for **Grocery & Gourmet Food : Wine : Red**



Renwood Winter Reds Port, Syrah,
Primitivo Mixed Pack, 3 x 750 mL

\$63.91 \$79.89

Eligible for 1¢ Standard Shipping [See Details](#)

[Show only Renwood items](#)

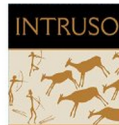


Renwood Delectable Port, Ice Wine,
Syrah Mixed Pack, 2 x 750 mL 1 x 375
mL

\$67.88 \$84.85

Eligible for 1¢ Standard Shipping [See Details](#)

[Show only Renwood items](#)



samsung android

mobile phones
features: android, samsung
samsung galaxy ace s5830 white
samsung galaxy 5 i5500 black
samsung galaxy 3 i5800
samsung galaxy s i9000 8gb ceramic white
samsung galaxy s i9000 8gb metallic black
samsung galaxy 5 i5500 white

tablet
features: android, samsung
samsung galaxy tab gt-p1000 16gb white

* translated from Italian

demo at <http://autocomplete.shoppydoo.com/demo.html>

[Porrini et al. WIAS 2014] R. Porrini, M. Palmonari and G. Vizzari. Composite Match Autocompletion (COMMA): a Semantic Result-Oriented Autocompletion Technique for e-Marketplaces. In *Web Intelligence and Agent Systems Journal*, 2014

[Palmonari et al. WI 2012] M. Palmonari, G. Vizzari, R. Porrini, A. Broglia, N. Lamberti. Comma: A Result-Oriented Composite Autocompletion Method for e-Marketplaces. In *Web Intelligence*, 2012

Multiple Classifications in Action - Product Autocomplete

samsung android|

mobile phones

features: android, samsung

samsung galaxy ace s5830 white

samsung galaxy 5 i5500 black

samsung galaxy 3 i5800

samsung galaxy s i9000 8gb ceramic white

samsung galaxy s i9000 8gb metallic black

samsung galaxy 5 i5500 white

tablet

features: android, samsung

samsung galaxy tab gt-p1000 16gb white

support for explorative keyword based queries

* translated from Italian

demo at <http://autocomplete.shoppydoo.com/demo.html>

[Porrini et al. WIAS 2014] R. Porrini, M. Palmonari and G. Vizzari. Composite Match Autocompletion (COMMA): a Semantic Result-Oriented Autocompletion Technique for e-Marketplaces. In *Web Intelligence and Agent Systems Journal*, 2014

[Palmonari et al. WI 2012] M. Palmonari, G. Vizzari, R. Porrini, A. Broglia, N. Lamberti. Comma: A Result-Oriented Composite Autocompletion Method for e-Marketplaces. In *Web Intelligence*, 2012

Multiple Classifications in Action - Product Autocomplete

samsung android

mobile phones

features: android, samsung

samsung galaxy ace s5830 white

samsung galaxy 5 i5500 black

samsung galaxy 3 i5800

samsung galaxy s i9000 8gb ceramic white

samsung galaxy s i9000 8gb metallic black

samsung galaxy 5 i5500 white

tablet

features: android, samsung

samsung galaxy tab gt-p1000 16gb white

* translated from Italian

support for explorative keyword based queries

facets and **categories** are considered in result-driven completion of the query



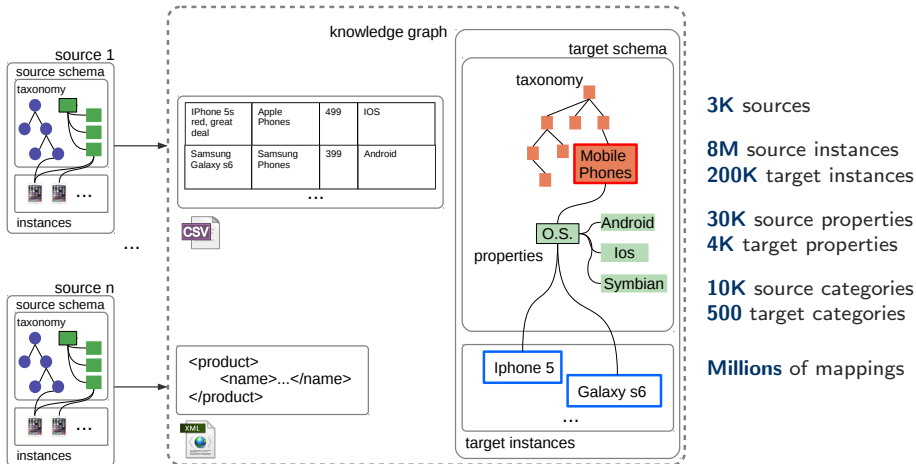
demo at <http://autocomplete.shoppydoo.com/demo.html>

[Porrini et al. WIAS 2014] R. Porrini, M. Palmonari and G. Vizzari. Composite Match Autocompletion (COMMA): a Semantic Result-Oriented Autocompletion Technique for e-Marketplaces. In *Web Intelligence and Agent Systems Journal*, 2014

[Palmonari et al. WI 2012] M. Palmonari, G. Vizzari, R. Porrini, A. Broglia, N. Lamberti. Comma: A Result-Oriented Composite Autocompletion Method for e-Marketplaces. In *Web Intelligence*, 2012

real-world example from eCommerce

Dataspace Management System



3K sources

8M source instances

200K target instances

30K source properties

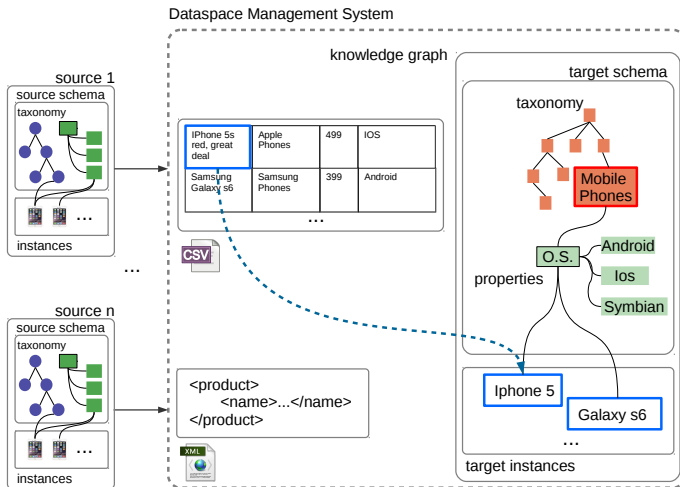
4K target properties

10K source categories

500 target categories

Millions of mappings

real-world example from eCommerce



3K sources

8M source instances

200K target instances

30K source properties

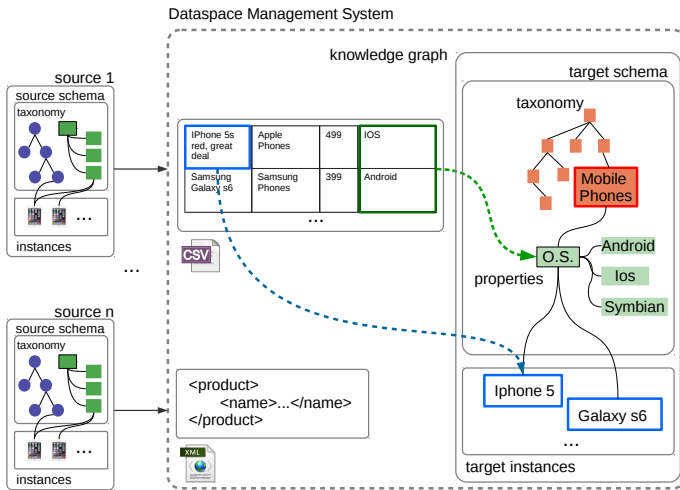
4K target properties

10K source categories

500 target categories

Millions of mappings

real-world example from eCommerce



3K sources

8M source instances

200K target instances

30K source properties

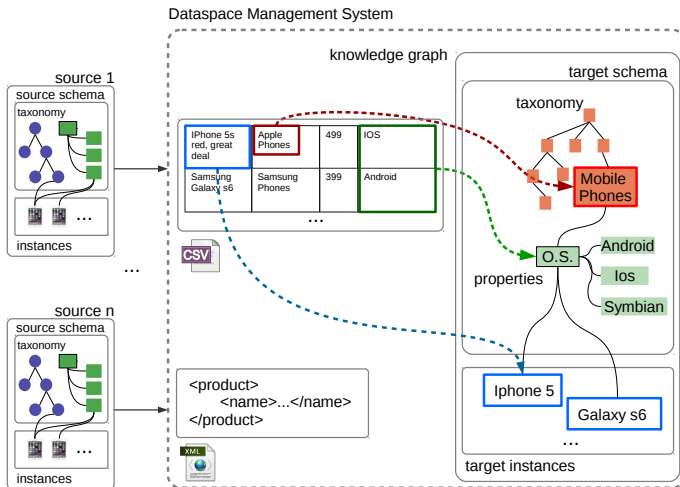
4K target properties

10K source categories

500 target categories

Millions of mappings

real-world example from eCommerce



3K sources

8M source instances

200K target instances

30K source properties

4K target properties

10K source categories

500 target categories

Millions of mappings

target schema and mapping maintenance is hard

- ▶ need for extensive knowledge of disparate domains
wines and clothes
- ▶ need for domain experts supervision
quality issues

target schema and mapping maintenance is hard

- ▶ need for extensive knowledge of disparate domains
wines and clothes
- ▶ need for domain experts supervision
quality issues

crucial for

- ▶ inclusion of new data sources covering different domains
- ▶ pay-as-you-go refinement of the integration
from wines
to italian, cabernet wine bottles, from 2011

highlighted by seminal works on Dataspaces [Franklin et al. 2005]

relevant to this thesis

- ▶ **schema enrichment**

- ▶ extraction of properties and categories

[Pound et al. 2011, Medelyan et al. 2013, Kong and Allan 2013] . . .

- ▶ category and property profiling

[Presutti et al. 2011, Jarrar et al. 2012] . . .

- ▶ **mapping discovery**

- ▶ schema and ontology matching (p-to-p, c-to-c)

[Bernstein et al. 2011, Shvaiko and Euzenat 2013] . . .

- ▶ web table annotation (sources with weak structure)

[Limaye et al. 2010, Venetis et al. 2011] . . .

Dataspace $\Delta = \langle \mathcal{S}, T, M \rangle$

- ▶ \mathcal{S} set of sources S with respective schema and instances
- ▶ T target knowledge graph (schema and instances)
- ▶ M set of mappings between the sources and the target

Dataspace $\Delta = \langle \mathcal{S}, T, M \rangle$

- ▶ \mathcal{S} set of sources S with respective schema and instances
- ▶ T target knowledge graph (schema and instances)
- ▶ M set of mappings between the sources and the target

schema	$\mathcal{A}^S = \langle C^S, P^S \rangle$	$\mathcal{A}^T = \langle C^T, P^T \rangle$
instances	$X^S = \{x_1^S, \dots, x_h^S\}$	$X^T = \{x_1^T, \dots, x_k^T\}$
categories	$C^S = \{c_1^S, \dots, c_m^S\}$	$C^T = \{c_1^T, \dots, c_n^T\}$
properties	$P^S = \{p_1^S, \dots, p_w^S\}$	$P^T = \{p_1^T, \dots, p_j^T\}$

- ▶ categories as FOL **unary** predicates from taxonomy or a lattice
MobilePhones(iphone 6)
- ▶ properties $p_i \subseteq X \times V_i$ as FOL **binary** predicates
price(iphone 6, "32 euro")

Dataspace $\Delta = \langle \mathcal{S}, T, M \rangle$

- ▶ \mathcal{S} set of sources S with respective schema and instances
- ▶ T target knowledge graph (schema and instances)
- ▶ M set of mappings between the sources and the target

Horn clauses encoded **mappings** from sources to target schema

$$a^T(\bar{x}) \leftarrow a_1^S(\bar{x}), \dots, a_i^S(\bar{x}), c_1, \dots, c_n$$

$a^T \in \mathcal{A}^T, a^S \in \mathcal{A}^S$
 c_1, \dots, c_n constraints

c-to-c (category-to-category)

$$Phones(x) \leftarrow ApplePhones(x)$$

p-to-p (property-to-property)

$$year(x, v) \leftarrow yearOfProduction(x, v)$$

c-to-p (category-to-property)

$$year(x, "2013") \leftarrow WinesFrom2013(x)$$

schema enrichment

- ▶ extraction of domain specific properties from sources

[Porrini et al. CAiSE 2014]

- ▶ analysis of property usage within the dataspace

[Palmonari et al. ESWC 2015]

schema enrichment

- ▶ extraction of domain specific properties from sources
[Porrini et al. CAiSE 2014]
- ▶ analysis of property usage within the dataspace
[Palmonari et al. ESWC 2015]

mapping discovery

- ▶ establishment of category-to-property mappings
[Porrini et al. CAiSE 2014]
- ▶ establishment of property-to-property mappings
Paper to be submitted

schema enrichment

- ▶ extraction of domain specific properties from sources
[Porrini et al. CAiSE 2014]
- ▶ analysis of property usage within the dataspace
[Palmonari et al. ESWC 2015]

mapping discovery

- ▶ establishment of category-to-property mappings
[Porrini et al. CAiSE 2014]
- ▶ establishment of property-to-property mappings
Paper to be submitted

case study from the eCommerce domain

- ▶ usage of target categories and properties for product autocompletion
[Palmonari et al. WI 2012, Porrini et al. WIAS 2014]

goal

granular characterization of dataspace instances by extracting domain specific properties

[Porrini et al. CAiSE 2014] R. Porrini, M. Palmonari and C. Batini. [Extracting Facets from Lost Fine-Grained Classifications in Dataspaces](#). In *CAiSE*, 2014

goal

granular characterization of dataspace instances by extracting domain specific properties

Wines

[Porrini et al. CAiSE 2014] R. Porrini, M. Palmonari and C. Batini. [Extracting Facets from Lost Fine-Grained Classifications in Dataspaces](#). In *CAiSE*, 2014

goal

granular characterization of dataspace instances by extracting domain specific properties

Wines

Winery Country of Origin	Wine Alcohol By Volume	Grape Variety	Wine Bottle Volume
<input type="checkbox"/> USA	<input type="checkbox"/> Under 10%	<input type="checkbox"/> Blend - White	<input type="checkbox"/> 375 mL
<input type="checkbox"/> China	<input type="checkbox"/> 10% to 12%	<input type="checkbox"/> Blend - Other	<input type="checkbox"/> 500 mL
<input type="checkbox"/> Australia	<input type="checkbox"/> 12% to 14%	<input type="checkbox"/> Fruit	<input type="checkbox"/> 750 mL
<input type="checkbox"/> Italy	<input type="checkbox"/> 14% & Up	<input type="checkbox"/> Muscadine	
Specialty Wine Type	Wine Vintage	<input type="checkbox"/> Cabernet Sauvignon	
<input type="checkbox"/> Sustainable	<input type="checkbox"/> 2011	<input type="checkbox"/> Pinot Noir	
<input type="checkbox"/> Small Lot	<input type="checkbox"/> 2010	<input type="checkbox"/> Chardonnay	
<input type="checkbox"/> Kosher	<input type="checkbox"/> 2009		
<input type="checkbox"/> Gluten-Free	<input type="checkbox"/> 2008		
	<input type="checkbox"/> 2007		

[Porrini et al. CAiSE 2014] R. Porrini, M. Palmonari and C. Batini. [Extracting Facets from Lost Fine-Grained Classifications in Dataspace](#). In *CAiSE*, 2014

goal

granular characterization of dataspace instances by extracting domain specific properties

observations about source categories

- ▶ source categories often come from specialized sources
(e.g., emarketplaces selling only wine bottles)
- ▶ c-to-c mappings typically map specialized to generic categories
 $Wines(x) \leftarrow Barolo(x)$

[Porrini et al. CAiSE 2014] R. Porrini, M. Palmonari and C. Batini. [Extracting Facets from Lost Fine-Grained Classifications in Dataspaces](#). In *CAiSE*, 2014

goal

granular characterization of dataspace instances by extracting domain specific properties

idea: leverage already defined c-to-c mappings

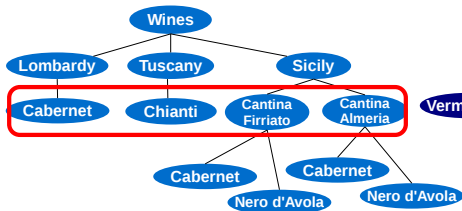
1. fix a target category c^T
2. pick all the source categories c^S such that $c^T(x) \leftarrow c^S(x)$
3. form clusters of **homogeneous** property values V^1, \dots, V^n
4. output properties $p_i \subseteq X \times V^i$
5. output c-to-p mappings: $p_i(x, v) \leftarrow v(x), v \in V^i$

[Porrini et al. CAISE 2014] R. Porrini, M. Palmonari and C. Batini. [Extracting Facets from Lost Fine-Grained Classifications in Dataspaces](#). In *CAISE*, 2014

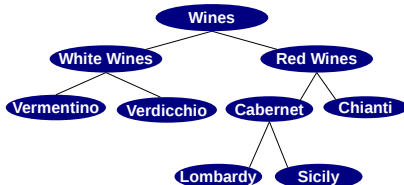
Source Category Mutual Exclusivity Principle

the more two source categories are mutually exclusive, the more they should be clustered together into the same property value set

Taxonomy of source S1



Taxonomy of source S2

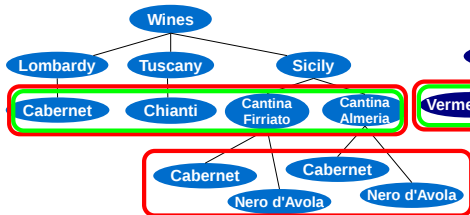


given two source categories c_1 and c_2 , their occurrence as siblings indicates that c_1 and c_2 are mutually exclusive

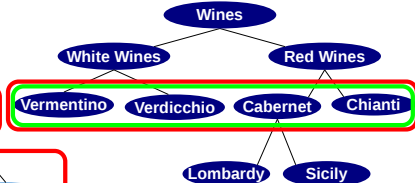
$$\text{TLD}(c_1, c_2) = 1 - \frac{|L_{c_1} \cap L_{c_2}|}{|L_{c_1} \cup L_{c_2}|}$$

Jaccard distance between the two sets of taxonomy layers where two categories c_1 and c_2 occur

Taxonomy of source S1



Taxonomy of source S2



cabernet *chianti*

$$\text{TLD}(\text{cabernet}, \text{chianti}) = 1 - \frac{|L_{\text{cabernet}} \cap L_{\text{chianti}}|}{|L_{\text{cabernet}} \cup L_{\text{chianti}}|} = 1 - \frac{2}{3} = \frac{1}{3}$$

	#Taxonomies	#Mappings
Wines	184	8967
Musical Instruments	128	1306
Grappe, Liquors, Aperitives	115	1254
Beers	58	156
DVD Movies	164	2042
Blu-Ray Movies	55	395
Dogs and Cats Food	80	5592
Rings	138	936
Ski and Snowboards	55	790
Necklaces	148	1156
Overall	688	22594

evaluation using real world data from the Italian PCE TrovaPrezzi

- ▶ property values manually grouped by domain experts
- ▶ comparison with Leacock and Chodorow similarity

[Leacock and Chodorow 1998]

- ▶ and with Wu and Palmer similarity

[Wu and Palmer 1994]

	<i>Value Effectiveness</i>			<i>Clustering Effectiveness</i>				<i>Quality</i>
	<i>P</i>	<i>R</i>	<i>F₁</i>	<i>F*</i>	<i>NMI*</i>	Purity	<i>E*</i>	<i>PRF*</i>
LC	0.394	0.953	0.537	0.666	0.709	0.220	0.685	0.531
WP	0.377	0.984	0.525	0.682	0.714	0.210	0.744	0.520
TLD	0.416	0.901	0.541	0.719	0.746	0.286	0.416	0.558

	<i>Value Effectiveness</i>			<i>Clustering Effectiveness</i>				<i>Quality</i>
	<i>P</i>	<i>R</i>	<i>F₁</i>	<i>F*</i>	<i>NMI*</i>	Purity	<i>E*</i>	<i>PRF*</i>
LC	0.394	0.953	0.537	0.666	0.709	0.220	0.685	0.531
WP	0.377	0.984	0.525	0.682	0.714	0.210	0.744	0.520
TLD	0.416	0.901	0.541	0.719	0.746	0.286	0.416	0.558

- ▶ TLD more effective in finding relevant property values and discarding noisy ones (high F_1)

	<i>Value Effectiveness</i>			<i>Clustering Effectiveness</i>				<i>Quality</i>
	<i>P</i>	<i>R</i>	<i>F₁</i>	<i>F*</i>	<i>NMI*</i>	Purity	<i>E*</i>	<i>PRF*</i>
LC	0.394	0.953	0.537	0.666	0.709	0.220	0.685	0.531
WP	0.377	0.984	0.525	0.682	0.714	0.210	0.744	0.520
TLD	0.416	0.901	0.541	0.719	0.746	0.286	0.416	0.558

- ▶ TLD more effective in finding relevant property values and discarding noisy ones (high F_1)
- ▶ TLD more effective in clustering homogeneous values (high *clustering effectiveness*)

Not Publishing ▾ View

Accepted Facets

New group... +

Regione ▾ ✕

New facet... +

- ✕ abruzzo
- ✕ argentina
- ✕ basilicata
- ✕ calabria
- ✕ campania
- ✕ emilia romagna
- ✕ francia
- ✕ friuli
- ✕ friuli venezia giulia
- ✕ lazio
- ✕ liguria
- ✕ lombardia
- ✕ marche
- ✕ piemonte
- ✕ puglia
- ✕ sardeana

Suggested Facets

 ↻

vino bianco +

vino rosato

vino rosso

bianco +

rosato

rosso

abruzzo +

argentina

basilicata

calabria

campania

emilia romagna

francia

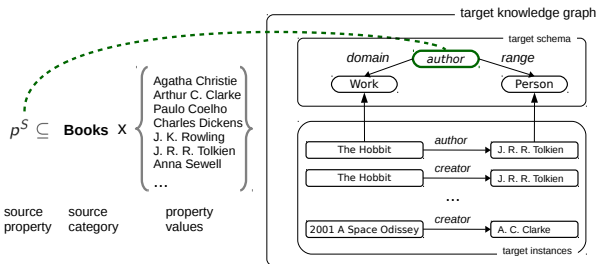
friuli

friuli venezia giulia

.

in production within the Italian PCE TrovaPrezzi dataspace
extracted properties for **16** target categories and counting

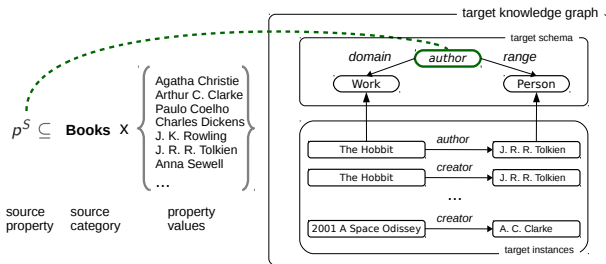
Property-to-Property Mapping



goal

given a source property p^S select a property p^T from the target schema such that the semantics of p^S is compatible with p^T

Property-to-Property Mapping



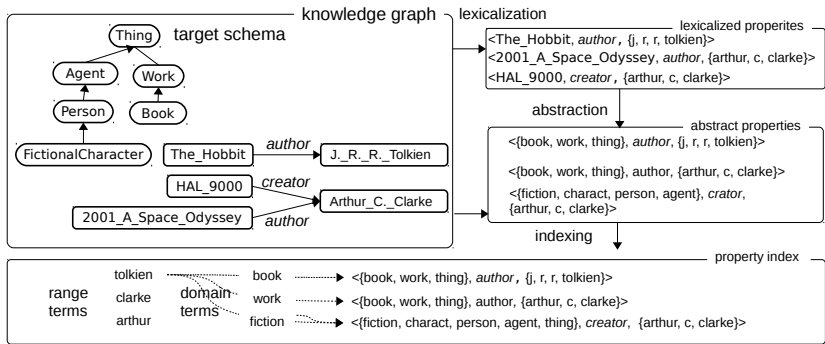
challenges

- ▶ more than one target property suitable for mapping
e.g., ~50000 properties from the DBpedia knowledge graph
- ▶ how to resolve ambiguities?

$$p^S \subseteq X^S \times \{2010, 2011, 2012, \dots\}$$

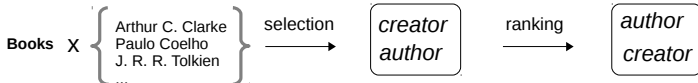
releaseYear for music albums
vintage for wines

- ▶ how to capture semantic compatibility between properties?



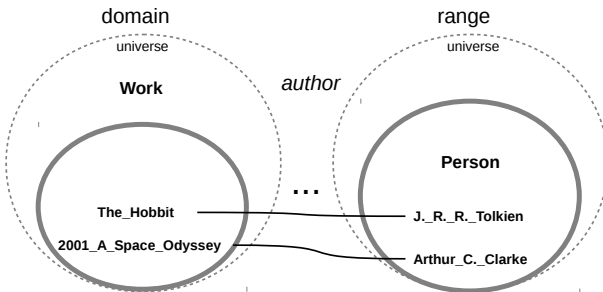
knowledge graph indexing

annotation



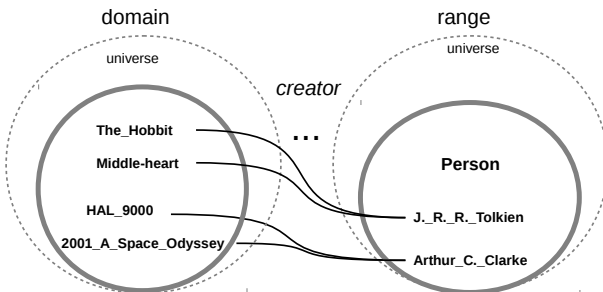
Specificity

$$p^S \subseteq \mathbf{Books} \times \left\{ \begin{array}{l} \text{Agatha Christie} \\ \text{Arthur C. Clarke} \\ \text{Paulo Coelho} \\ \text{Charles Dickens} \\ \text{J. K. Rowling} \\ \text{J. R. R. Tolkien} \\ \text{Anna Sewell} \\ \dots \end{array} \right\}$$



Specificity

$$p^S \subseteq \mathbf{Books} \times \left\{ \begin{array}{l} \text{Agatha Christie} \\ \text{Arthur C. Clarke} \\ \text{Paulo Coelho} \\ \text{Charles Dickens} \\ \text{J. K. Rowling} \\ \text{J. R. R. Tolkien} \\ \text{Anna Sewell} \\ \dots \end{array} \right\}$$



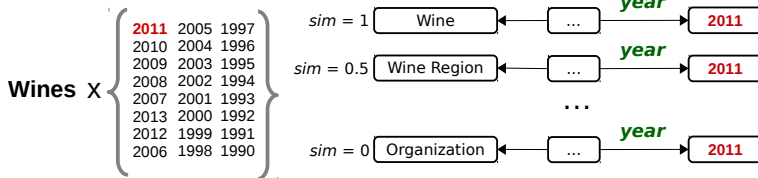
Coverage



Coverage



Frequency



	<i>source properties</i>		<i>target properties</i>		
	#	categories	#	relevant	fair
dbpedia-numbers	8	7	53195	~4	~19
dbpedia-entities	31	13	53195	~7	~7
dbpedia	39	13	53195	~3	~9
yago-full	83	17	89	1	-
yago-abstract	83	10	89	1	-

- ▶ DBPedia and YAGO as target knowledge graphs
- ▶ DBPedia based gold standard created through a questionnaire
- ▶ YAGO based state-of-the-art gold standard

[Limaye et al. 2010, Venetis et al. 2011]

Mean Average Precision on the DBPedia Gold Standard

	dbpedia-numbers	dbpedia-entities	dbpedia
majority	0.22	0.50	0.44
maximum likelihood	0.16	0.39	0.34
proposed approach	0.25	0.55*	0.49*

Mean Reciprocal Rank on the YAGO Gold Standard

	yago-full	yago-abstract
majority	0.76	0.86
maximum likelihood	0.81	0.85
proposed approach	0.88*	0.90*

* $p < 0.05$

- ▶ comparison with the well known majority voting scheme and
- ▶ with maximum likelihood based approach [Venetis et al. 2011]

STAN Alpha

Semantic Table Annotation Tool

[Home](#) / CPS_Schools_2013-2014_Academic_Year.csv

Annotation [Ontologies](#) [Namespace](#)

Save

Export ▾

School	name	FullName	SchoolName2	ISBE Name	address	Street Direction	address	city	S
400010	Ace Technical Chtr HS	Architecture, Construction, and Engineering(ACE)Technical Charter School	Ace Technical Chtr HS	Ace Technical Charter High School	5410	S	State St	Chicago	IL
609772	Addams	Jane Addams Elementary School	Addams	Addams Elem School	10810	S	Avenue H	Chicago	IL
609773	Agassiz	Louis A Agassiz Elementary School	Agassiz	Agassiz Elem School	2851	N	Seminary Ave	Chicago	IL
610513	Air Force HS	Air Force Academy High School	Air Force HS	Air Force Acad High School	3630	S	Wells St	Chicago	IL
610212	Albany Park	Albany Park Multicultural Academy	Albany Park	Albany Park Multicultural Elem	4929	N	Sawyer Ave	Chicago	IL
609774	Alcott ES	Louisa May Alcott	Alcott ES	Alcott Elem	2625	N	Orchard St	Chicago	IL

demo @ <http://stan.disco.unimib.it>

algorithm source code @ <http://bitbucket.org/rporrini/cluster-labelling>

credits to Brando Preda (MSc student) for the web app development

goal

support the analysis of property usage in a knowledge graph

- ▶ what categories are described in the knowledge graph?
- ▶ what properties are used to characterize the instances?
- ▶ how frequent is the use of a given property?

[Palmonari et al. ESWC 2015] M. Palmonari, A. Rula, R. Porrini, A. Maurino, B. Spahiu and V. Ferme. **ASBTAT**: Linked Data Summaries with ABstraction and STATistics. In *ESWC Posters and Demos*, 2015

ABSTAT linked data summaries



browse the abstract knowledge patterns found into

subject type (occurrences)		predicate (occurrences)	object type (occurrences)	frequency
<input type="button" value="filter"/>	<input type="text" value="subject"/>	<input type="text" value="predicate"/>	<input type="text" value="object"/>	
	http://www.ontologydesignpatterns.org/ont/dul/DUL.owl#Agent (1688266)	http://xmlns.com/foaf/0.1/name (3546622)	http://www.w3.org/2000/01/rdf-schema#Literal (10407268)	2270174
	http://xmlns.com/foaf/0.1/Person (1445775)	http://xmlns.com/foaf/0.1/name (3546622)	http://www.w3.org/2000/01/rdf-schema#Literal (10407268)	2058043
	http://schema.org/Person (1445106)	http://xmlns.com/foaf/0.1/name (3546622)	http://www.w3.org/2000/01/rdf-schema#Literal (10407268)	2057209
	http://wikidata.dbpedia.org/resource/Q215627 (1445106)	http://xmlns.com/foaf/0.1/name (3546622)	http://www.w3.org/2000/01/rdf-schema#Literal (10407268)	2057209
	http://wikidata.dbpedia.org/resource/Q5 (1445106)	http://xmlns.com/foaf/0.1/name (3546622)	http://www.w3.org/2000/01/rdf-schema#Literal (10407268)	2057209
	http://www.ontologydesignpatterns.org/ont/dul/DUL.owl#NaturalPerson (1445106)	http://xmlns.com/foaf/0.1/name (3546622)	http://www.w3.org/2000/01/rdf-schema#Literal (10407268)	1980978
	http://xmlns.com/foaf/0.1/Person (1445775)	http://purl.org/dc/elements/1.1/description (1738736)	http://www.w3.org/2000/01/rdf-schema#Literal (10407268)	1738440
	http://www.ontologydesignpatterns.org/ont/dul/DUL.owl#Agent (1688266)	http://purl.org/dc/elements/1.1/description (1738736)	http://www.w3.org/2000/01/rdf-schema#Literal (10407268)	1738356
	http://wikidata.dbpedia.org/resource/Q215627 (1445106)	http://purl.org/dc/elements/1.1/description (1738736)	http://www.w3.org/2000/01/rdf-schema#Literal (10407268)	1737966
	http://schema.org/Person (1445106)	http://purl.org/dc/elements/1.1/description (1738736)	http://www.w3.org/2000/01/rdf-schema#Literal (10407268)	1737966

demo @ <http://abstat.disco.unimib.it>

source code @ <http://github.com/rporrini/abstat>

studied the problem of **property management** in dataspace

- ▶ extraction of domain specific properties
- ▶ establishment of **c-to-p** and **p-to-p** mappings
- ▶ evaluation on different application domains

eCommerce - LOD

- ▶ algorithms deployed in production or in research prototypes

last mile(s)

- ▶ refine the formal framework of the thesis
- ▶ experiment **p-to-p** mapping approach in the eCommerce domain

Publications

- [1] M. Palmonari, A. Rula, R. Porrini, A. Maurino, B. Spahiu and V. Ferme. **ASBTAT: Linked Data Summaries with ABstraction and STATistics**. In *ESWC*, 2015
- [2] R. Porrini, M. Palmonari and C. Batini. **Extracting Facets from Lost Fine-Grained Classifications in Dataspace**. In *CAiSE*, 2014
- [3] R. Porrini, M. Palmonari and G. Vizzari. **Composite Match Autocompletion (COMMA): a Semantic Result-Oriented Autocompletion Technique for e-Marketplaces**. In *Web Intelligence and Agent Systems Journal*, 2014
- [4] M. Palmonari, G. Vizzari, R. Porrini, A. Broglia, N. Lamberti. **Comma: A Result-Oriented Composite Autocompletion Method for e-Marketplaces**. In *Web Intelligence*, 2012

Conference Paper Reviews (as sub-reviewer)

- ▶ **ISWC 2015**: 15th International Semantic Web Conference
- ▶ **ESWC 2015**: 12th Extended Semantic Web Conference
- ▶ **EDBT 2015**: 18th International Conference on Extending Database Technology
- ▶ **AAAI-15**: 23th AAAI Conference on Artificial Intelligence
- ▶ **EKAW 2014**: 19th International Conference on Knowledge Engineering and Knowledge Management
- ▶ **ISWC 2014**: 14th International Semantic Web Conference
- ▶ **WI 2014**: 2014 IEEE/WIC/ACM International Conference on Web Intelligence
- ▶ **AAAI-14**: 22th AAAI Conference on Artificial Intelligence
- ▶ **ESWC 2014**: 11th Extended Semantic Web Conference
- ▶ **CAISE 2014**: 26th International Conference on Advanced Information Systems Engineering
- ▶ **ODBASE 2013**: 12th International Conference on Ontologies, DataBases, and Applications of Semantics
- ▶ **WI 2013**: 2013 IEEE/WIC/ACM International Conference on Web Intelligence

Courses and Schools

- ▶ The Impact of Logic: from Proof Systems to Databases - Politecnico di Milano (*evaluation pending*)
- ▶ Recommender Systems - DISCo - (*final report submitted*)
- ▶ Cluster Analysis - DISCo
- ▶ Third ESWC Summer School - Kalamaky - Crete (GR)
- ▶ Foundations of Data Exchange and Integration - Politecnico di Milano
- ▶ Advanced Analytics and Behavior Informatics - DISCo

Seminars

- ▶ From Sentiment Analysis to Continuous Learning - DISCo
- ▶ Introduzione alla Logica nella Rappresentazione della Conoscenza - DISCo
- ▶ La Misurazione Della Felicità al Tempo dei Big Data - DISCo
- ▶ Optimum Hyperpaths in Directed Hypergraphs - DISCo
- ▶ Phase Transitions in Social and Economic Systems - DISCo
- ▶ Progettare e Fare Open Data. Metodologia e tools sviluppati in Evodevo a partire dall'esperienza Open Data INPS - DISCo
- ▶ Semantic Constraints for Data Quality Assessment and Cleaning - DISCo
- ▶ Big Data e la forza degli eventi - DISCo
- ▶ WOA 2012: 13th National Workshop "Dagli Oggetti agli Agenti" - DISCo

Teaching

- ▶ Tutor for *Data and Web Semantics* course (MSc and PhD, Fall 2014) - University of Illinois at Chicago
- ▶ Co-Advisor of two BSc thesis and two MSc thesis
- ▶ Lecturer (2 lectures) for the "*Artificial Intelligence*" course (MSc A.A. 2012/2013 and 2013/2014) - DISCo
- ▶ Tutor for *Distributed Systems* course (BSc A.A. 2012/2013) - DISCo

Questions?

Source taxonomies are:

- ▶ **many**
3900 within the TrovaPrezzi italian price comparison engine
- ▶ **noisy**
type > white > by vine > chardonnay > producer > firriato
- ▶ **heterogeneous**
type > white > by vine > chardonnay > producer > firriato
wines > white wines > greco di tufo
- ▶ **ambiguous** different semantics for different contexts
red is a wine type for wines
and a color for shirts

- ▶ *document corpora*

focus on property hierarchies - specific for unstructured data
[Stoica et al. 2007, Dakka and Ipeirotis 2008, Wei et al. 2013, Medelyan et al. 2013]

- ▶ search engines' *query logs* and *documents*

user search queries as a primary source of information
[Li et al. 2009, Pasca and Alfonseca 2009, Pound et al. 2011]

- ▶ search engines' *query results*

integrate and rank properties already present in web documents
[Yan et al. 2010, Dou et al. 2011, Kawano et al. 2012, Kong and Allan 2013]

Similarity-Relatedness between taxonomy categories

- ▶ Leacock and Chodorow similarity [Leacock and Chodorow 1998]

- ▶ Wu and Palmer similarity [Wu and Palmer 1994]

- ▶ ...

not designed for taxonomies

from ontologies

- ▶ intensional matchers - focus on schema, neglect instances
[Cheatham and Hitzler 2014] . . .
- ▶ extensional matchers - focus on the instances, neglect the schema
[Zhang et al. 2015] . . .

from tabular data

- ▶ custom knowlege graphs - knowledge graphs with special features
[Venetis et al. 2011, Wang et al. 2012] . . .
- ▶ holistic approaches - annotate the table as a whole
[Limaye et al. 2010, Mulwad et al. 2013, Zhang 2015] . . .

General References

- [Franklin et al. 2005] M. Franklin, A. Halevy and D. Maier. From Databases to Dataspaces: A New Abstraction for Information Management. In *SIGMOD Record*, 2005
- [Bernstein et al. 2011] P. A. Bernstein, J. Madhavan and E. Rahm. Generic Schema Matching, Ten Years Later. In *PVLDB*, 2011
- [Presutti et al. 2011] V. Presutti, L. Aroyo, A. Adamou, B. A. C. Schopman, A. Gangemi and G. Schreiber. Extracting Core Knowledge from Linked Data. In *COLD*, 2011
- [Jarrar et al. 2012] M. Jarrar and M.D. Dikaiakos. A Query Formulation Language for the Data Web. In *IEEE Trans. Knowl. Data Eng.*, 2012
- [Shvaiko and Euzenat 2013] Pavel Shvaiko and Jérôme Euzenat. Ontology Matching: State of the Art and Future Challenges. In *IEEE Trans. Knowl. Data Eng.*, 2013

Property Extraction References

- [Wu and Palmer 1994] Z. Wu and M. Palmer. Verb semantics and lexical selection. In *ACL*, 1994
- [Leacock and Chodorow 1998] C. Leacock and M. Chodorow. Combining local context and wordnet similarity for word sense identification. In *MIT Press*, 1998
- [Stoica et al. 2007] E. Stoica, M.A. Hearst and M. Richardson. Automating creation of hierarchical faceted metadata structures. In *HLT-NAACL*, 2007
- [Dakka and Ipeirotis 2008] W. Dakka and P.G. Ipeirotis. Automatic extraction of useful facet hierarchies from text databases. In *ICDE*, 2008
- [Li et al. 2009] X. Li, Y.Y. Wang, A. Acero. Extracting structured information from user queries with semi-supervised conditional random fields. In *SIGIR*, 2009
- [Pasca and Alfonseca 2009] M. Pasca and E. Alfonseca. Web-derived resources for web information retrieval: from conceptual hierarchies to attribute hierarchies. In *SIGIR*, 2009

- [Yan et al. 2010] N. Yan, C. Li, S.B. Roy, R. Ramegowda and G. Das. *Facetedpedia: enabling query-dependent faceted search for wikipedia*. In *CIKM*, 2010
- [Pound et al. 2011] J. Pound, S. Paparizos and P. Tsaparas. *Facet discovery for structured web search: a query-log mining approach*. In *SIGMOD*, 2011
- [Dou et al. 2011] Z. Dou, S. Hu, Y. Luo, R. Song and J.R. Wen. *Finding dimensions for queries*. In *CIKM*, 2011
- [Kawano et al. 2012] Y. Kawano, H. Ohshima and K. Tanaka. *On-the-fly generation of facets as navigation signs for web objects*. In *DASFAA*, 2012
- [Wei et al. 2013] B. Wei, J. Liu, J. Ma, Q. Zheng, W. Zhang and B. Feng. *Dft-extractor: a system to extract domain-specific faceted taxonomies from wikipedia*. In *WWW*, 2013
- [Medelyan et al. 2013] O. Medelyan, S. Manion, J. Broekstra, A. Divoli, A.L. Huang and I.H. Witten. *Constructing a focused taxonomy from a document collection*. In *ESWC*, 2013
- [Kong and Allan 2013] W. Kong and J. Allan. *Extracting query facets from search results*. In *SIGIR*, 2013

Property Mapping References

- [Limaye et al. 2010] G. Limaye, S. Sarawagy and S. Chakrabarti. *Annotating and Searching Web Tables Using Entities, Types and Relationship*. In *VLDB*, 2010
- [Venetis et al. 2011] P. Venetis, A. Halevy, J. Madhavan, M. Pasca, W. Shen, F Wu, G. Miao and C. Wu. *Recovering Semantics of Tables on the Web*. In *VLDB*, 2011
- [Wang et al. 2012] J. Wang, H. Wang, Z. Wang, and K. Q. Zhu. *Understanding tables on the web*. In *ER*, 2012
- [Mulwad et al. 2013] V. Mulwad, T. Finin, and A. Joshi. *Semantic message passing for generating linked data from tables*. In *ISWC*, 2013
- [Cheatham and Hitzler 2014] M. Cheatham and P. Hitzler. *The properties of property alignment*. In *ISWC*, 2014
- [Zhang et al. 2015] Z. Zhang, A. L. Gentile, I. Augenstein, E. Blomqvist, and F. Ciravegna. *An Unsupervised Data-driven Method to Discover Equivalent Relations in Large Linked Datasets*. In *Sem. Web - accepted for publication*, 2015

[Zhang 2015] Z. Zhang. Effective and Efficient Semantic Table Interpretation using TableMiner+. In *Sem. Web - under transparent review*, 2015