

Aligning Unions of Concepts in Ontologies of Geospatial Linked Data

Rahul Parundekar, Craig A. Knoblock and Jose-Luis Ambite
{parundek,knoblock,ambite}@usc.edu
University of Southern California

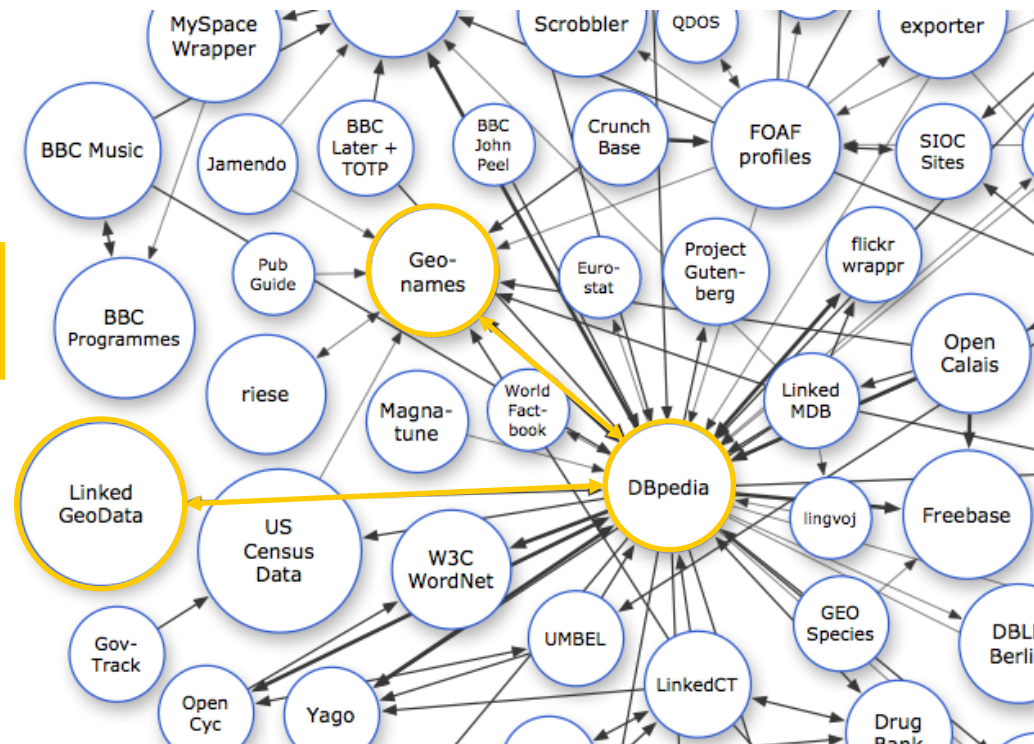
INTRODUCTION



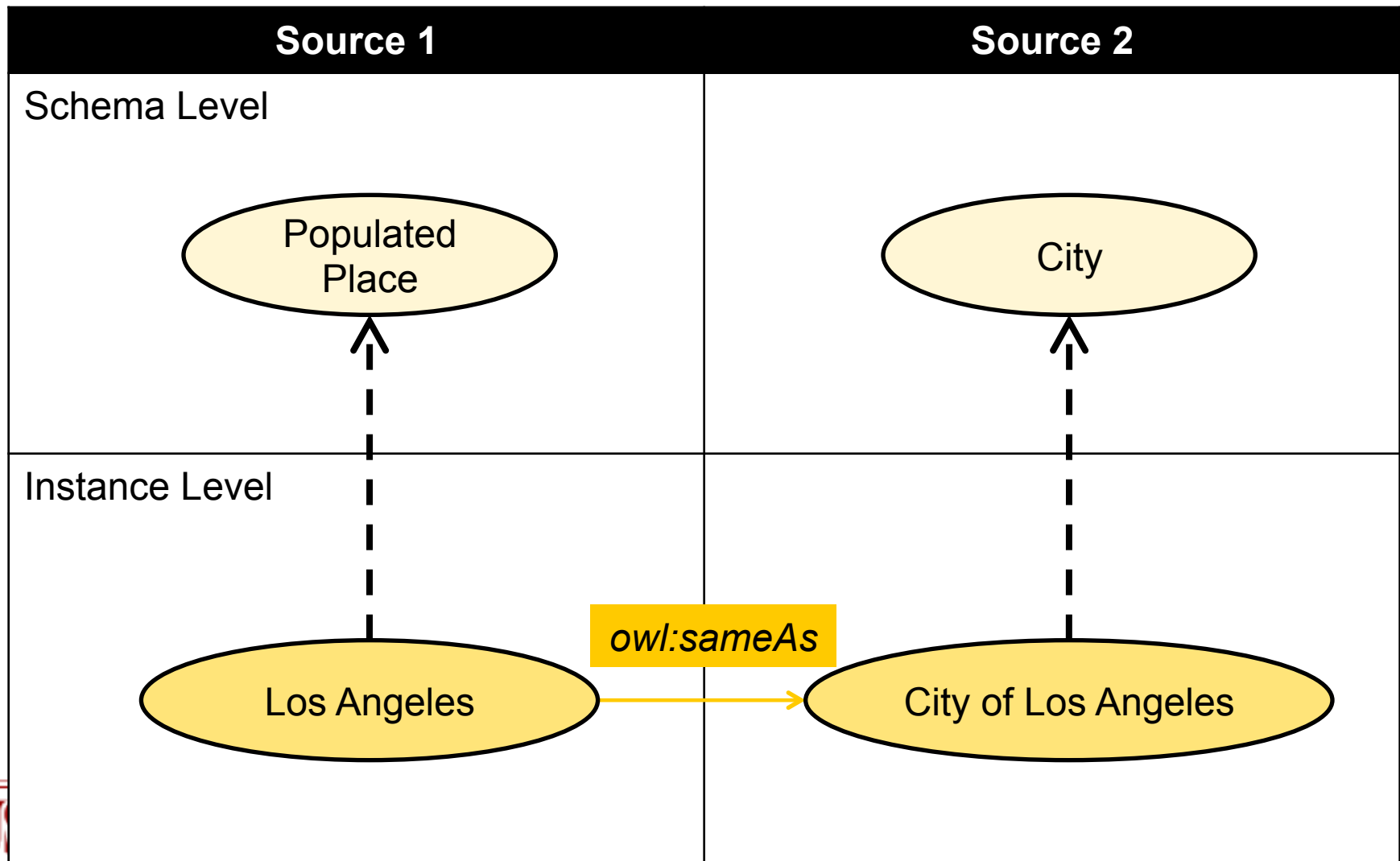
Web of Geospatial Linked Data

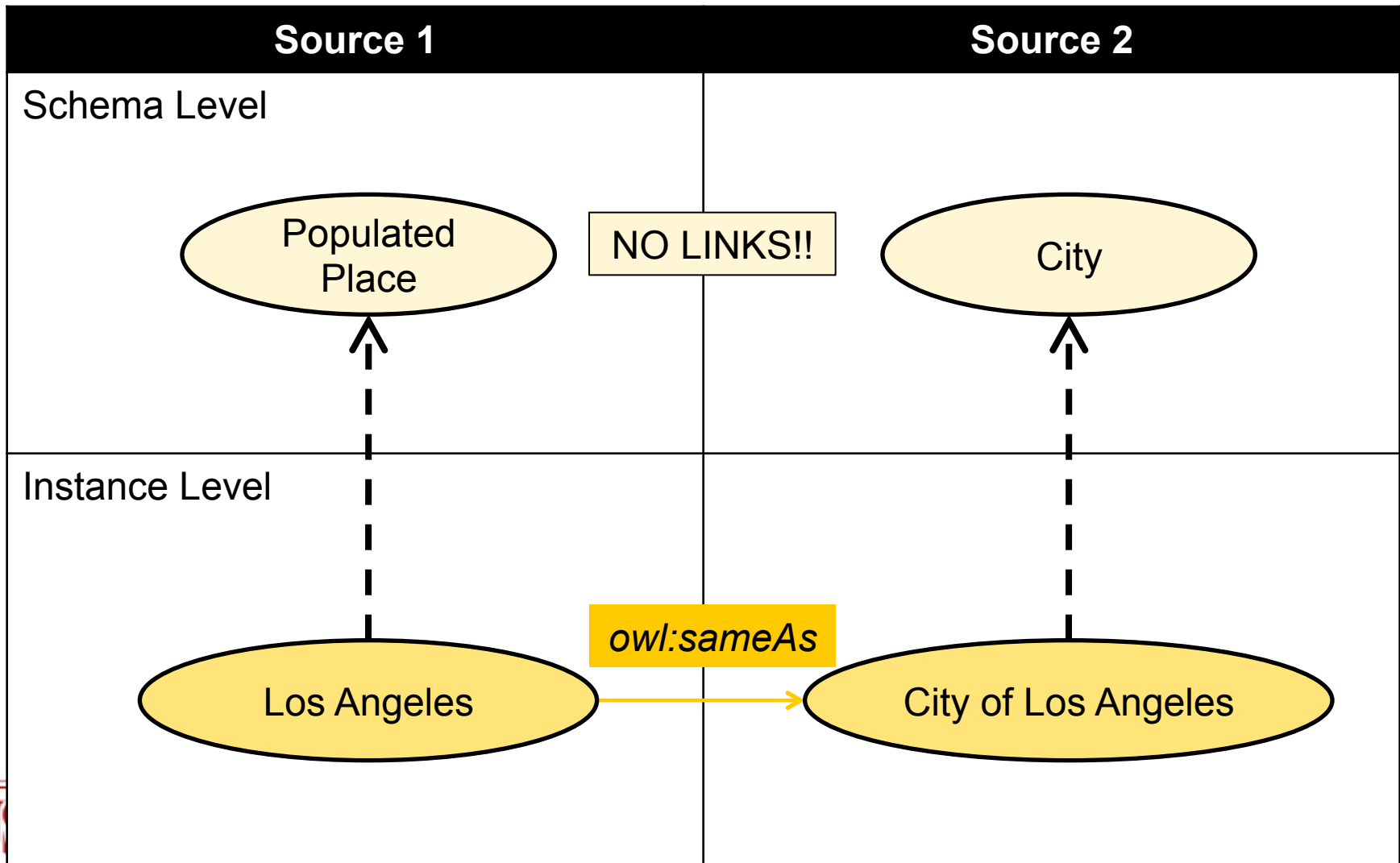
- Different sources with different schemas
- Equivalent instances in the geospatial domain connected with *owl:sameAs*

Geospatial
Domain

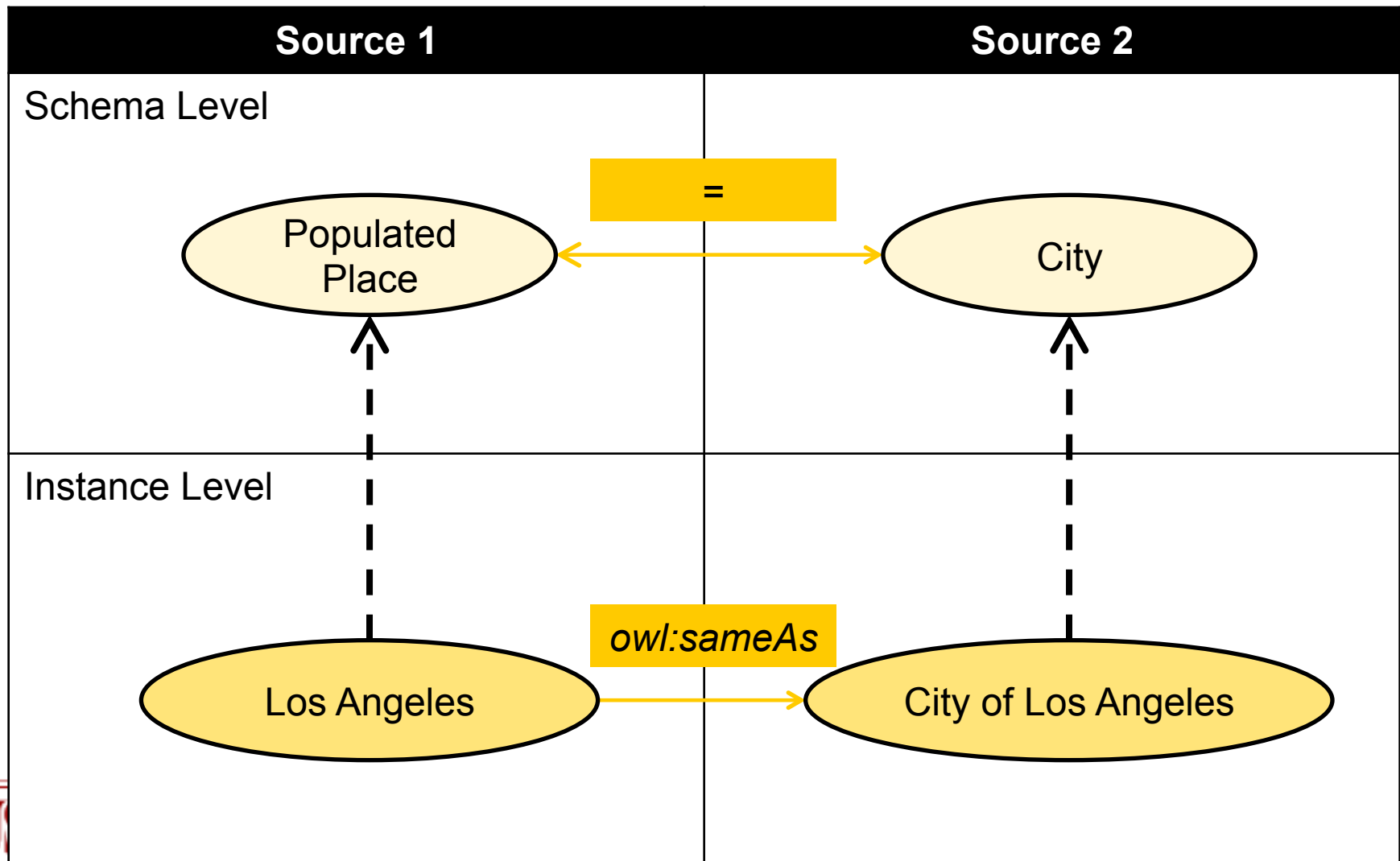


Interlinked instances...





Can we find schema alignments?



Previous Work @ ISWC 2010

Linking and Building Ontologies of Linked Data



Extensional Approach to Ontology Alignment

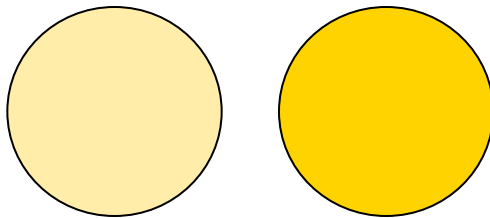


Represents set of instances belonging to ClassA

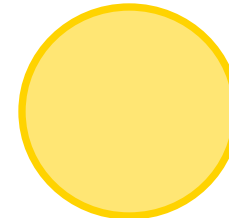


Represents set of instances belonging to ClassB

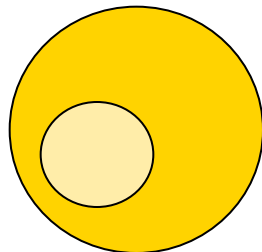
ClassA is disjoint from ClassB



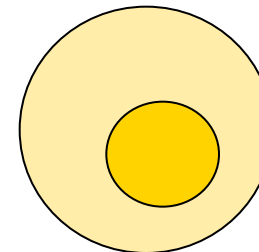
ClassA is equivalent to ClassB



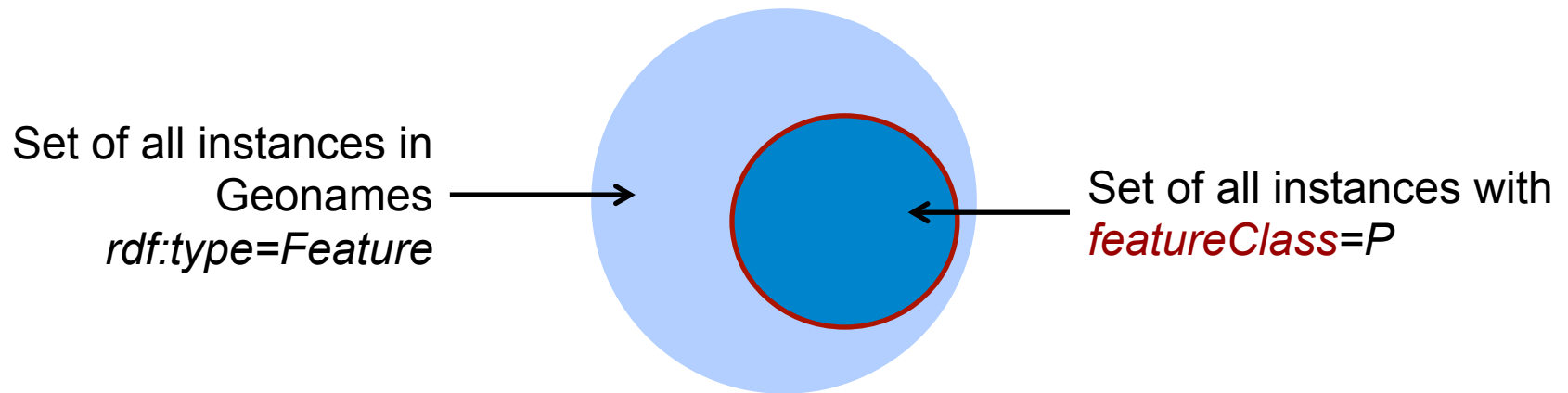
ClassA is subset of ClassB



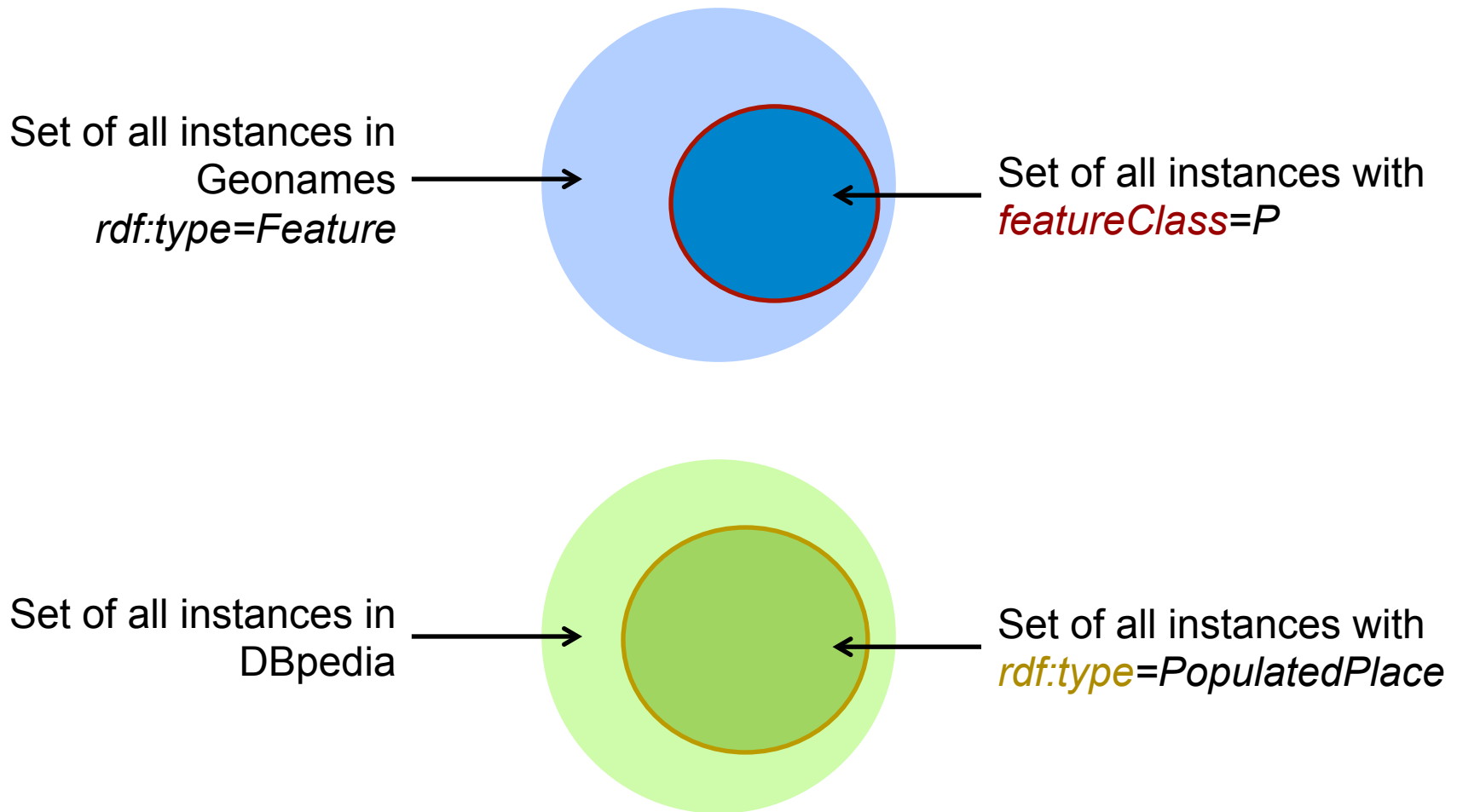
ClassB is subset of ClassA



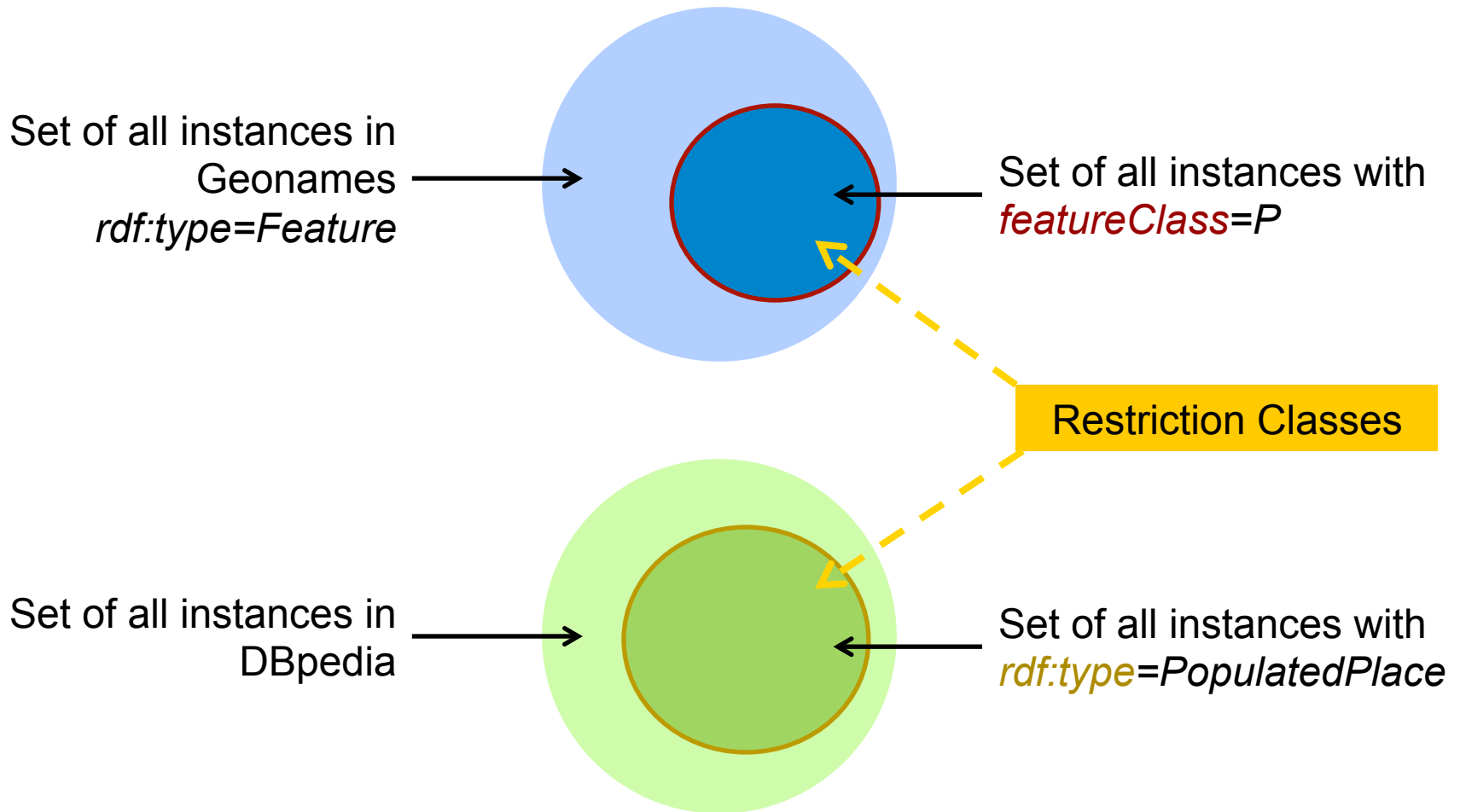
Classes are created extensionally by adding value restrictions on properties



Classes are created extensionally by adding value restrictions on properties

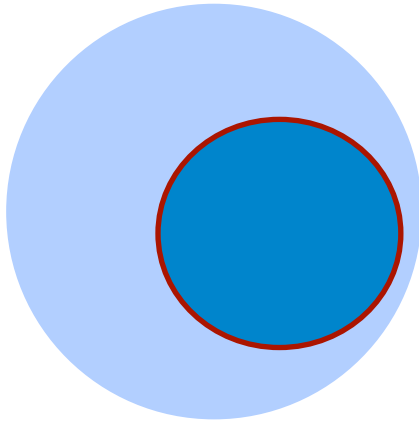


Classes are created extensionally by adding value restrictions on properties



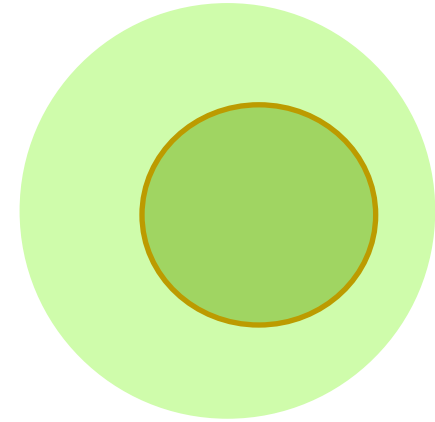
Aligning Restriction Classes Using Extensional Approach

featureClass=P



r_1

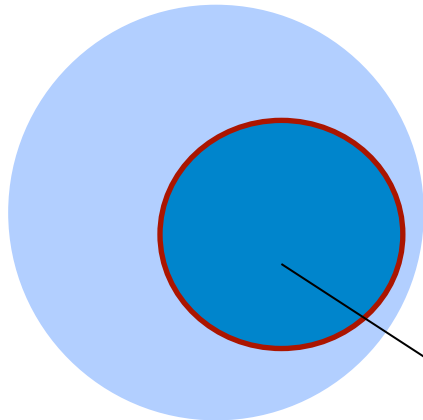
rdf:type=PopulatedPlace



r_2

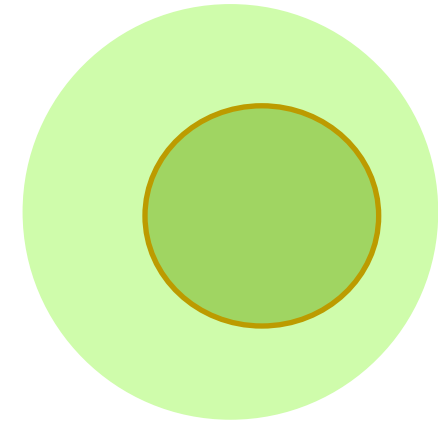
Aligning Restriction Classes Using Extensional Approach

featureClass=P



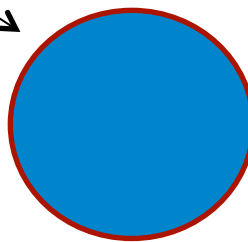
r_1

rdf:type=PopulatedPlace



r_2

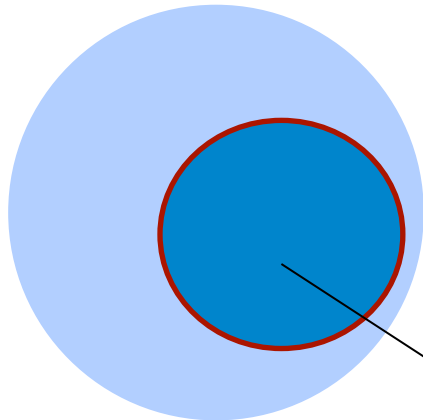
$\text{Img}(r_1)$



Set of instances from DBpedia
that r_1 is linked to

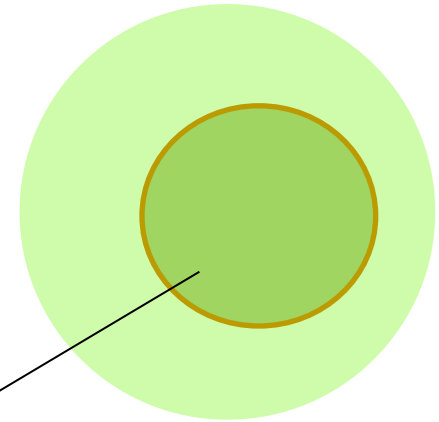
Aligning Restriction Classes Using Extensional Approach

featureClass=P

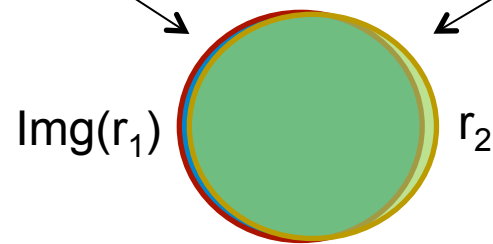


r_1

rdf:type=PopulatedPlace



r_2



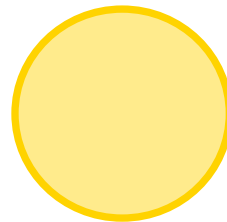
Extensionally, when are two classes equal?



Represents set of instances belonging to ClassA



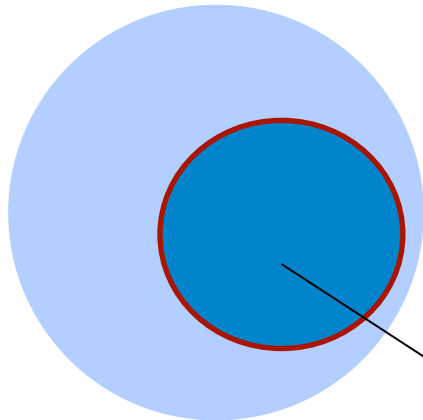
Represents set of instances belonging to ClassB



$$\frac{|\text{ClassA} \cap \text{ClassB}|}{|\text{ClassA}|} = \frac{|\text{ClassA} \cap \text{ClassB}|}{|\text{ClassB}|} = 1$$

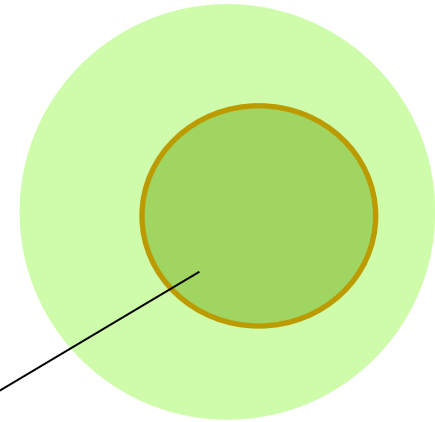
Aligning Restriction Classes Using Extensional Approach

featureClass=P

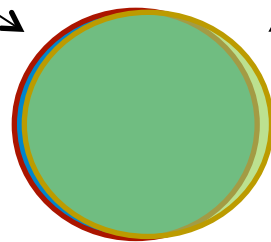


r_1

rdf:type=PopulatedPlace



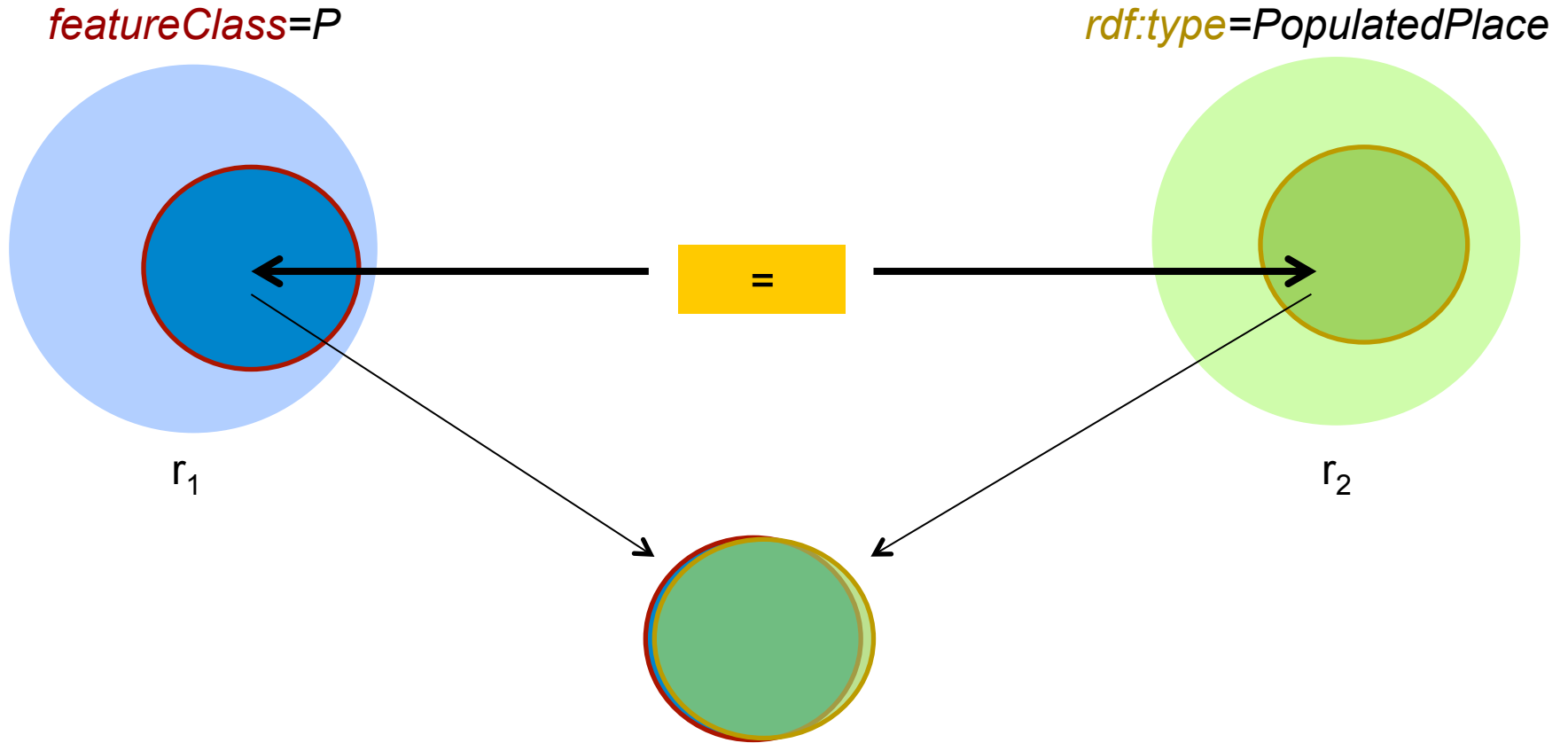
r_2



$$\frac{|\text{Img}(r_1) \cap r_2|}{|\text{Img}(r_1)|} > 0.9$$

$$\frac{|\text{Img}(r_1) \cap r_2|}{|r_2|} > 0.9$$

Aligning Restriction Classes Using Extensional Approach



$$\frac{|\text{Img}(r_1) \cap r_2|}{|\text{Img}(r_1)|} > 0.9$$

$$\frac{|\text{Img}(r_1) \cap r_2|}{|r_2|} > 0.9$$

- Algorithm was able to
 - Specialize ontologies where original were rudimentary
 - Find complimentary hierarchy across an ontology
- Alignments based on the actual data
 - reflects the semantics of the sources in practice
- Equivalences, Subset alignments before and after removing implied alignments

Source 1 (O_1)	Source 2 (O_2)	$\#(r_1 = r_2)$ total	$\#(r_1 = r_2)$ best matches	$\#(r_1 \subset r_2)$ before	$\#(r_1 \subset r_2)$ after	$\#(r_2 \subset r_1)$ before	$\#(r_2 \subset r_1)$ after
LinkedGeoData	DBpedia	158	152	2528	1837	1804	1627
Geonames	DBpedia	31	19	809	400	1384	1247
Geospecies	DBpedia	509	420	9112	2294	6098	4455
MGI	GeneID	10	9	2031	1869	3594	2070
Geospecies	Geospecies	94	88	1550	1201	-	-

- Algorithm was able to
 - Specialize ontologies where original were rudimentary
 - Find complimentary hierarchy across an ontology
- Alignments based on the actual data
 - reflects the semantics of the sources in practice
- Equivalences, Subset alignments before and after removing implied alignments

Source 1 (O_1)	Source 2 (O_2)	$\#(r_1 = r_2)$ total	$\#(r_1 = r_2)$ best matches	$\#(r_1 \subset r_2)$ before	$\#(r_1 \subset r_2)$ after	$\#(r_2 \subset r_1)$ before	$\#(r_2 \subset r_1)$ after
LinkedGeoData	DBpedia	158	152	2528	1837	1804	1627
Geonames	DBpedia	31	19	809	400	1384	1247
Geospecies	DBpedia	509	420	9112	2294	6098	4455
MGI	GeneID	10	9	2031	1869	3594	2070
Geospecies	Geospecies	94	88	1550	1201	-	-

Can we use the subset relations to find more meaningful alignments?



TerraCognita Workshop - ISWC 2011

Aligning Unions of Concepts in Ontologies of Linked Geospatial Data



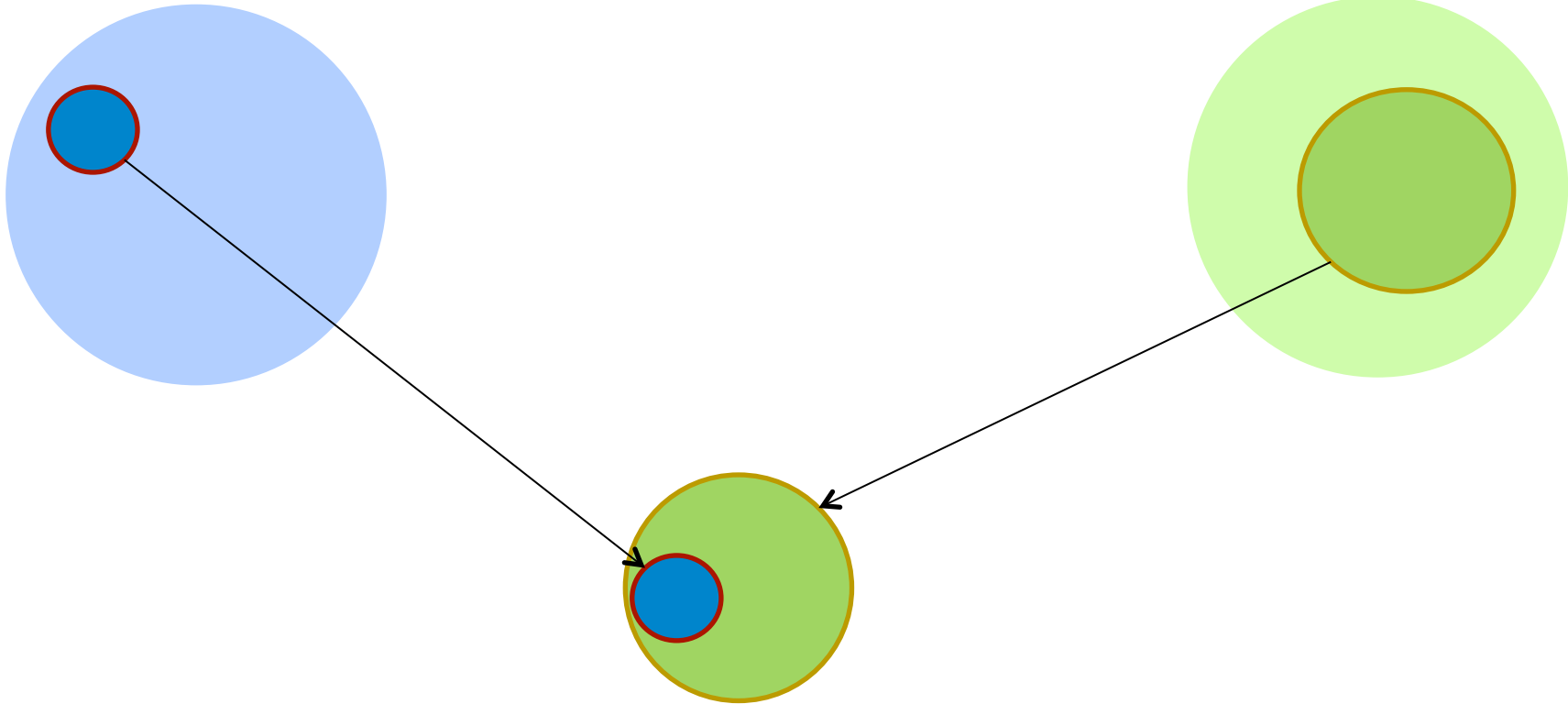
Is there a pattern in the subset relations?

Let's look at 3 of the subset relations we found...

1) Schools in *Geonames* are Educational Institutions in *DBpedia*

featureCode=S.SCH

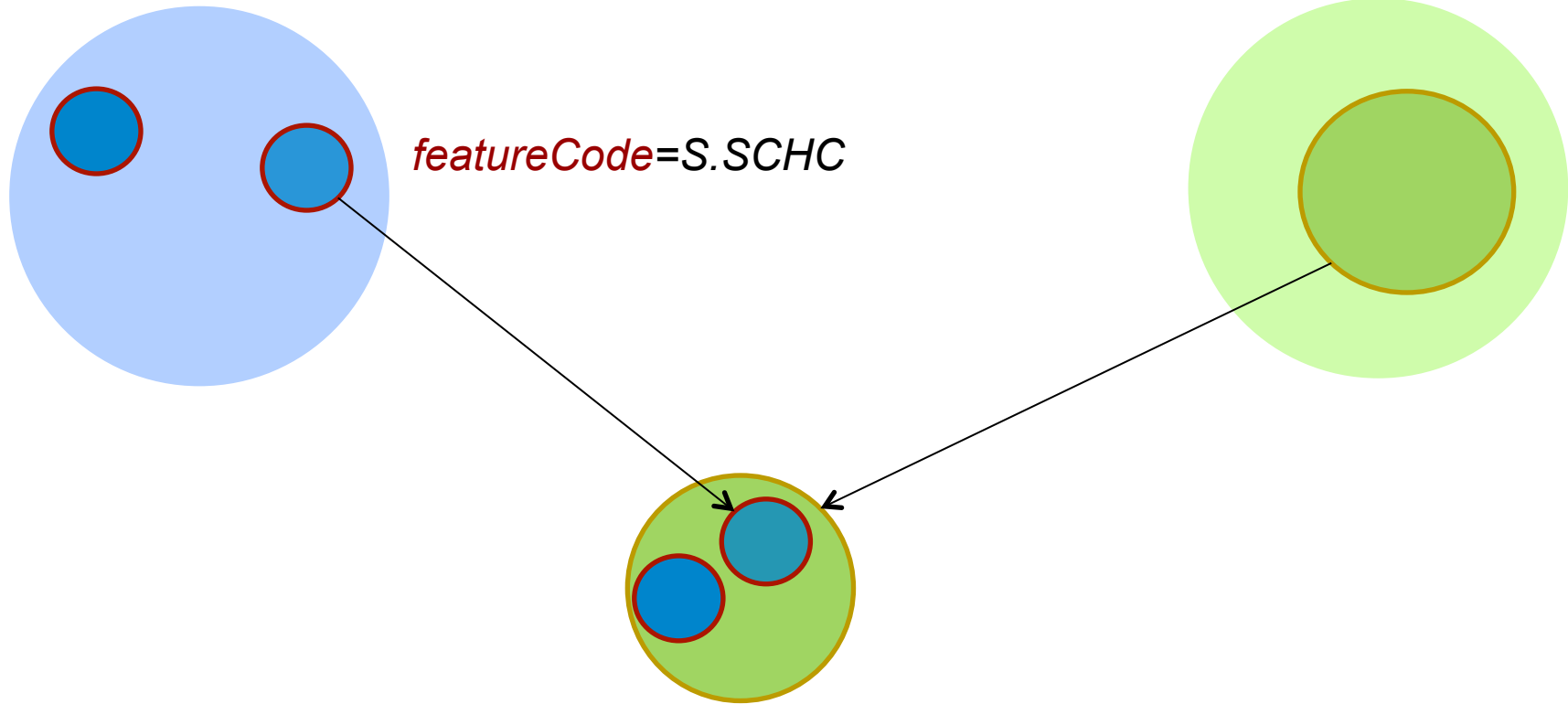
rdf:type=EducationalInstitution



2) Colleges in *Geonames* are Educational Institutions in *DBpedia*

featureCode=S.SCH

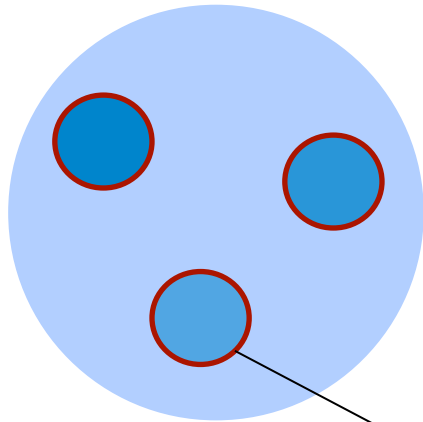
rdf:type=EducationalInstitution



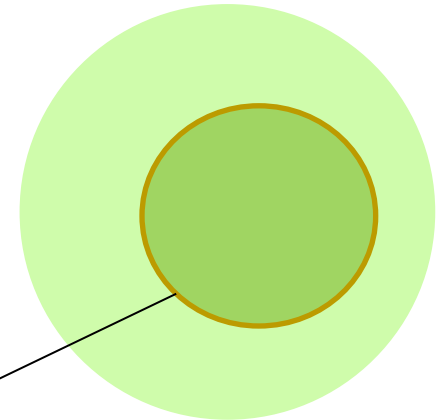
3) Universities in *Geonames* are Educational Institutions in *DBpedia*

featureCode=S.SCH

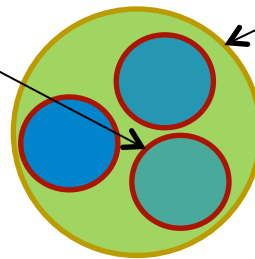
rdf:type=EducationalInstitution



featureCode=S.SCHC



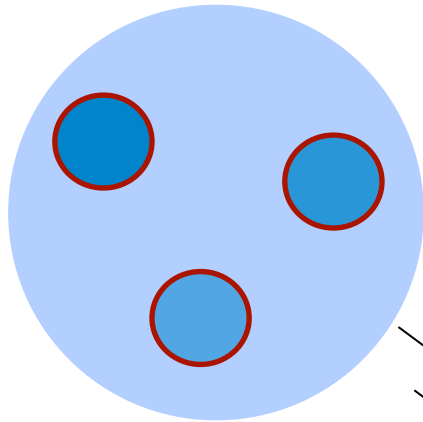
featureCode=S.UNIV



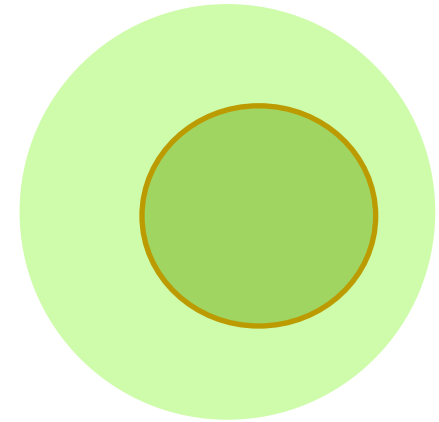
Taken by themselves, the subset relations are not useful

featureCode=S.SCH

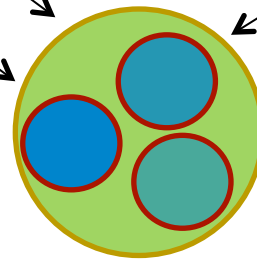
rdf:type=EducationalInstitution



featureCode=S.SCHC

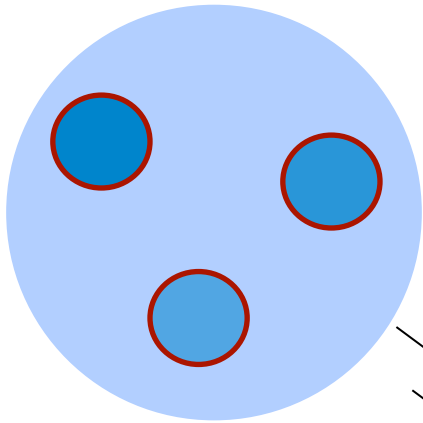


featureCode=S.UNIV



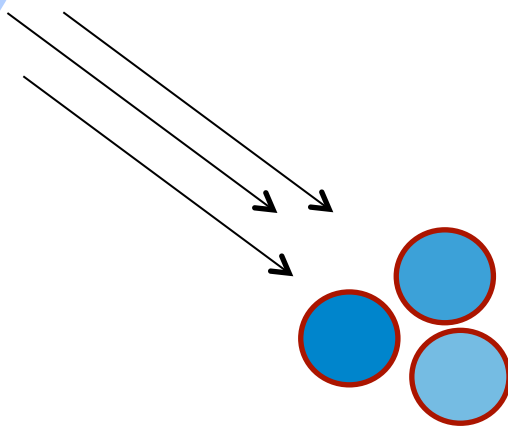
We use the common *featureCode* property as a hint...

featureCode=S.SCH

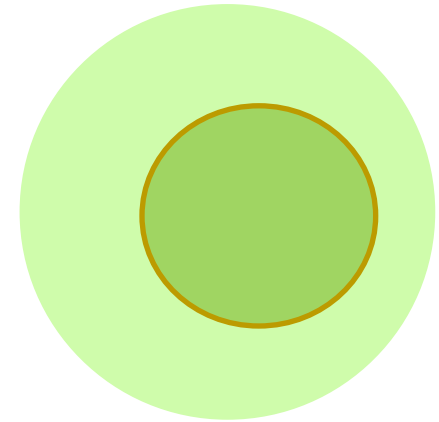


featureCode=S.SCHC

featureCode=S.UNIV



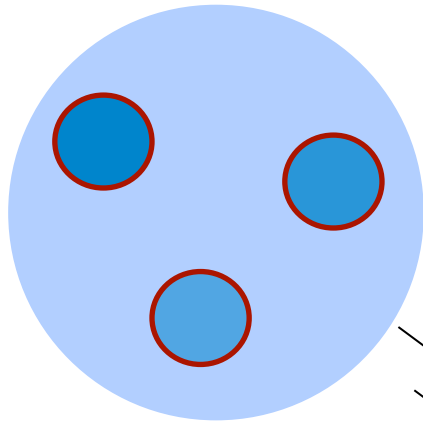
rdf:type=EducationalInstitution



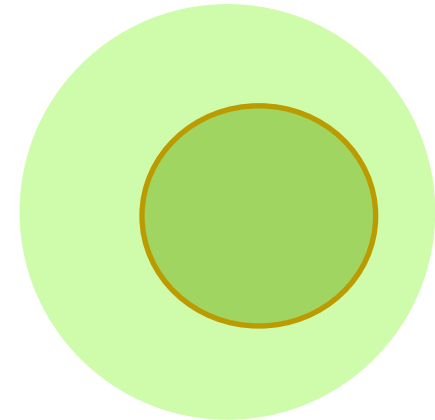
...to form a *Union* of Restriction Classes

featureCode=S.SCH

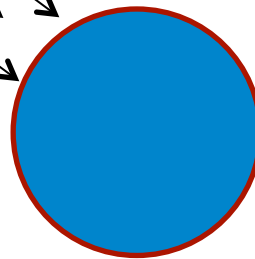
rdf:type=EducationalInstitution



featureCode=S.SCHC



featureCode=S.UNIV

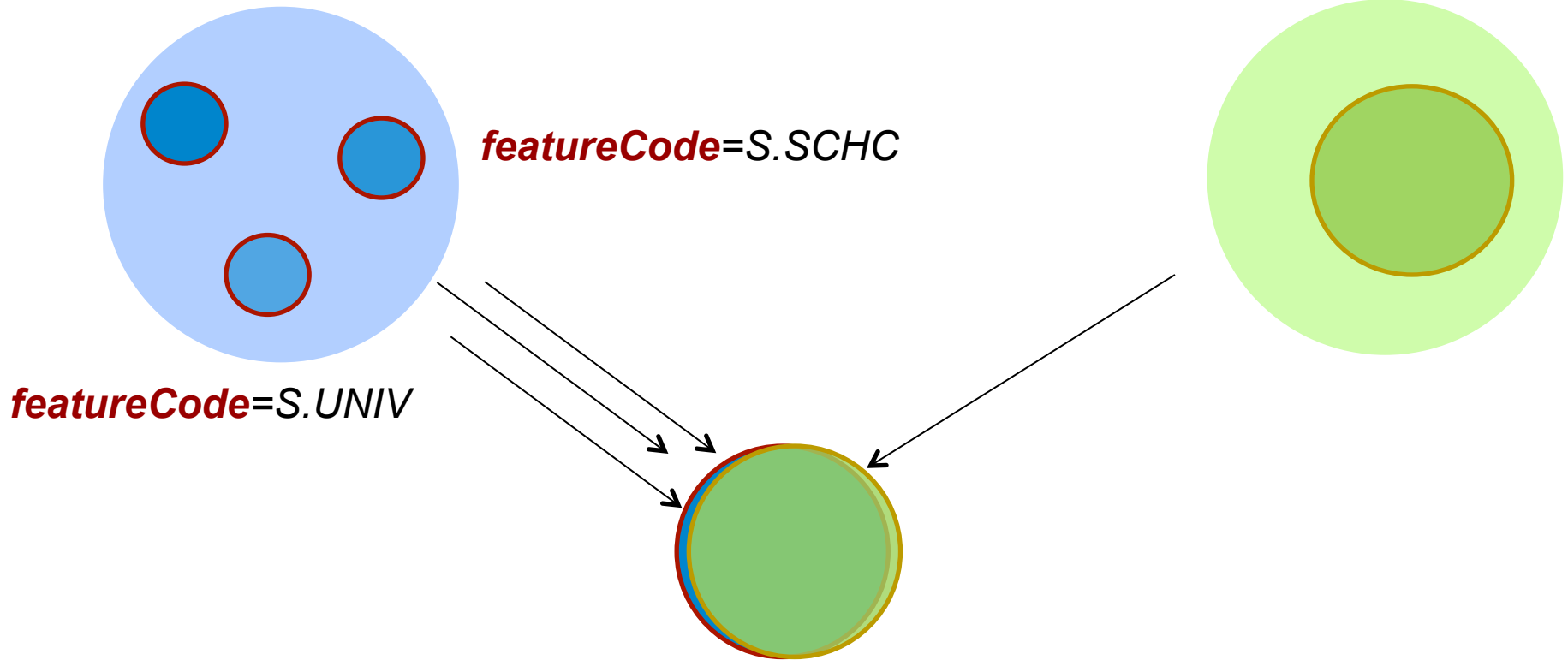


featureCode=S.SCH **U** *featureCode*=S.SCHC **U** *featureCode*=S.UNIV

Contribution 1: Find Union Alignments

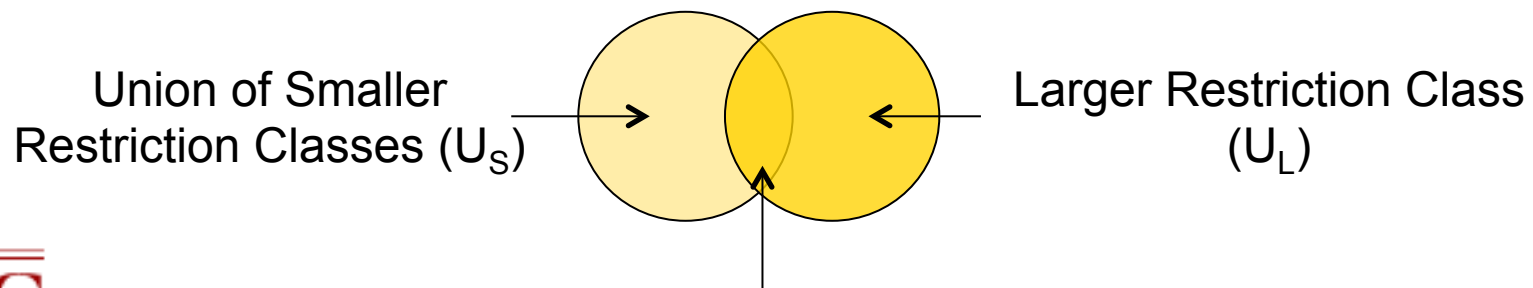
featureCode=S.SCH

rdf:type=EducationalInstitution



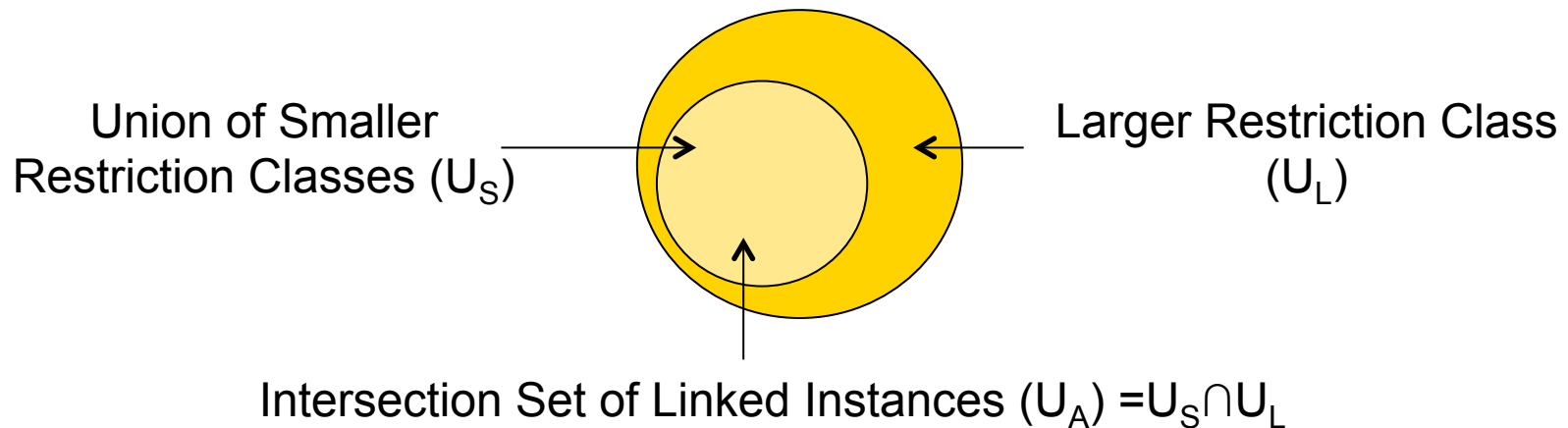
featureCode=S.SCH \cup *featureCode*=S.SCHC \cup *featureCode*=S.UNIV

- For all alignments found in the ISWC2010 paper marked as subsets
 1. We group all subset alignments according to the common larger restriction class
 2. We form a *union concept* such that all restriction classes
 - have the same property
 - have a single *property-value pair* each
 3. We then try to match the *union concept* to the larger class
 4. This forms a hypothesis *Union Alignment*



Intersection Set of Linked Instances (U_A) = $U_S \cap U_L$

Finding Union Alignments: Scoring



$$\frac{|U_A|}{|U_S|} = 1 \text{ since by definition, all smaller classes are subsets}$$

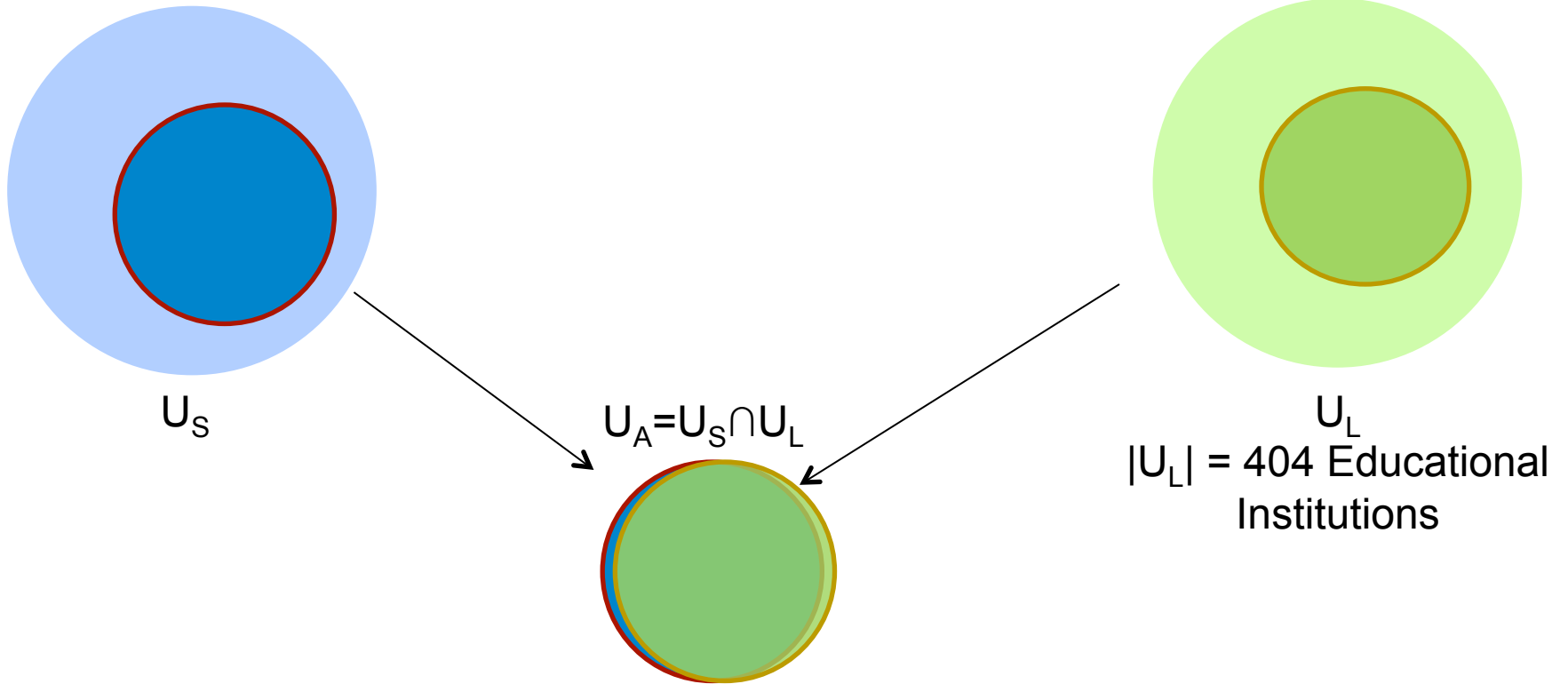
So, if $\frac{|U_A|}{|U_L|} = 1$, then the larger class U_L is equivalent to U_S

Practically, we use a relaxed subset assumption: $\frac{|U_A|}{|U_S|}, \frac{|U_A|}{|U_L|} > 0.9$

Contribution 1: Find Union Alignments

featureCode={S.SCH, S.SCHC, S.UNIV}

rdf:type=EducationalInstitution



$$\frac{|U_A|}{|U_S|} > 0.9$$

$$\frac{|U_A|}{|U_L|} = \frac{396}{404} = 0.98 > 0.9$$

What are the other 8 Educational Institutions?

Contribution 2: Find Outliers / Discrepancies

- We are also able to point out where the instances that disagree with the alignment lie
- These instances were not part of the alignment because
 - Their restriction class was not a subset ($P' < 0.9$)
 - Some of these instances are
 - Linked Incorrectly with *owl:sameAs*
 - Assigned wrong value during RDF generation*
 - Common in both sets (could be debatable)
 - Did not have a minimum support size of 2 instances (set with 1 instance cannot be relied on)
- Outliers help in understanding discrepancies in the Linked Data

What are the other 8 Educational Institutions?

- 1 with *featureCode*=S.HSP (Hostpitals)
 - There are 31 instances with S.HSP because of which Hospitals are not subsets
- 3 with *featureCode*=S.BLDG (Buildings)
- 1 with *featureCode*=S.EST (Establishment)
- 1 with *featureCode*=S.LIBR (Library)
- 1 with *featureCode*=S.MUS (Museum)
- 1 doesn't have a *featureCode* property

RESULTS



Larger class from *DBpedia* and union of smaller classes from *Geonames*

#	Sub-group $\{p_1, v_1, p_2\}$	List(v_2)	$R'_U = \frac{ U_A }{ U_L }$	$ U_A $	$ U_L $	Outliers	# Explained Instances
1	{ <i>rdf:type</i> , <i>dbpedia:EducationalInstitution</i> , <i>geonames:featureCode</i> }	S.SCH, S.SCHC, S.UNIV	0.9801	396	404	S.BLDG (3/122), S.EST (1/13), S.LIBR (1/7), S.HSP (1/31), S.MUS (1/43)	403
2	{ <i>dbpedia:country</i> , <i>dbpedia:Spain</i> , <i>geonames:countryCode</i> }	ES	0.9997	3917	3918	IT (1/7635)	3918
3	{ <i>dbpedia:region</i> , <i>dbpedia:Basse-Normandie</i> , <i>geonames:parentADM2</i> }	geonames:2989247, geonames:2996268, geonames:3029094	1.0	754	754		754
4	{ <i>rdf:type</i> , <i>dbpedia:Airport</i> , <i>geonames:featureCode</i> }	S.AIRB, S.AIRP	0.9924	1981	1996	S.AIRF (9/22), S.FRMT (1/5), S.SCH (1/404), S.STNB (2/5), S.STNM (1/36), T.HLL (1/61)	1996

Larger class from *Geonames* and union of smaller classes from *DBpedia*

#	Sub-group $\{p_1, v_1, p_2\}$	List(v_2)	$R'_U = \frac{ U_A }{ U_L }$	$ U_A $	$ U_L $	Outliers	# Explained Instances
5	{ <i>geonames:countryCode</i> , <i>NL</i> , <i>dbpedia:country</i> }	dbpedia:Netherlands, dbpedia:The_Netherlands, dbpedia:Flag_of_the _Netherlands.svg	0.9802	1939	1978	dbpedia:Kingdom_of _the_Netherlands	1940
6	{ <i>geonames:countryCode</i> , <i>JO</i> , <i>dbpedia:country</i> }	dbpedia:Jordan dbpedia:Flag_of_Jordan.svg	0.95	19	20		20

Larger class from *DBpedia* and union of smaller classes from *LinkedGeoData*

#	Sub-group $\{p_1, v_1, p_2\}$	List(v_2)	$R'_U = \frac{ U_A }{ U_L }$	$ U_A $	$ U_L $	Outliers	# Explained Instances
7	{ <i>dbpedia:bundesland</i> , <i>Saarland</i> , <i>lgd:OpenGeoDBLicensePlateNumber</i> }	HOM, IGB, MZG, NK, SB, SLS, VK, WND	0.93	46	49		46
8	{ <i>rdf:type</i> , <i>dbpedia:EducationalInstitution</i> , <i>rdf:type</i> }	<i>lgd:Amenity</i> , <i>lgd:K2543</i> , <i>lgd:School</i> , <i>lgd:University</i> , <i>lgd:WaterTower</i>	0.9901	2609	2610		2609

Larger class from *LinkedGeoData* and union of smaller classes from *DBpedia*

#	Sub-group $\{p_1, v_1, p_2\}$	List(v_2)	$R'_U = \frac{ U_A }{ U_L }$	$ U_A $	$ U_L $	Outliers	# Explained Instances
9	{ <i>lgd:gnisST_alpha</i> , <i>NJ</i> , <i>dbpedia:subdivisionName</i> }	Atlantic, Burlington, Cape May, Hudson, Hunterdon, Monmoth, New Jersey, Ocean, Passaic	1.0	214	214		214
10	{ <i>rdf:type</i> , <i>lgd:Waterway</i> , <i>rdf:type</i> }	<i>dbpedia:Stream</i> , <i>dbpedia:River</i>	0.97	33	34	<i>dbpedia:Place(1/94989)</i>	34

We find a total of 6595 Union Alignments

Source 1	Source 2	Larger Class from Source 1	Larger Class from Source 2	Total number of Union Alignments found
Geonames	DBpedia	434	318	752
LinkedGeoData	DBpedia	2746	3097	5843

Results also available at

<http://www.isi.edu/integration/data/UnionAlignments>

- **BLOOMS, BLOOMS+ ([4][5] in paper)**
 - Linked Open Data ontologies aligned with ‘Proton’
 - Constructs a forest of concepts and computes structural similarity
 - Geonames – Proton has “poor performance” because of small number and vague classes in Geonames
- **Volker et al. ([8] in paper)**
 - Statistical schema induction
 - Mines associativity rules from intermediate *‘transaction datasets’*
 - Develops OWL2 Axioms
- **AgreementMaker [2]**
 - Similarity Metrics on labels of classes

- **Conclusion**
 - We were able to find *Union Alignments* in the Geospatial Domain
 - Find alignments where no direct equivalence was evident
 - Introduced a disjunction operator to restriction classes
 - We were able to find *Outliers*
 - Help identify inconsistencies in the data
- **Future work**
 - Our algorithm is not limited to Geospatial domain. We would like to explore other domains
 - Experimental comparison with other approaches
 - Preliminary findings suggest patterns in properties like *geonames:countryCode* and *dbpedia:country*

Any questions?

THANK YOU

