



Entity Resolution in the Web of Data



Vassilis Christophides

christop@csd.uoc.gr

<http://www.csd.uoc.gr/~hy562>

University of Crete, Fall 2019



Describing and Linking Entities: Entity-Centric Applications & Knowledge Bases



 **Google** bilbao 

All Images Maps News Videos More ▾ Search tools

About 91,200,000 results (0.51 seconds)

Bilbao - Wikipedia, the free encyclopedia
<https://en.wikipedia.org/wiki/Bilbao> ▾
 Bilbao is a municipality and city in Spain, a major city in the province of Biscay in the autonomous community of the Basque Country. It is the largest municipality ...
 Guggenheim Museum Bilbao - Athletic Bilbao - Metro Bilbao - Bilbao la Vieja

Tourists - Bilbao Turismo
www.bilbaoturismo.net/BilbaoTurismo/en/tourists ▾
 Night falls, the city is waiting for you but you need to sleep...Everything from extreme luxury and comfort to very economical options is available in Bilbao.
 BILBAO IN THREE DAYS - How to get there and move ... - Leisure and Agenda

Images for bilbao 



More images for bilbao

 2

 **Google** bilbao points of interest 

The Top 10 Things to Do in Bilbao - TripAdvisor - Bilbao, Spain ...
https://www.tripadvisor.com/Attractions-g187454-Activities-Bilbao_Vizcaya_Provinc... ▾
 Hotels near Guggenheim Museum Bilbao. Hotels near Casco Viejo. Hotels near Plaza Nueva. Hotels near Funicular de Artxanda. Hotels near Bilbao Fine Arts Museum. Hotels near Azkuna Zentroa. Hotels near White Bridge (Zubi Zur) Hotels near La Ribera Market.
 Casco Viejo - Guggenheim Museum Bilbao - Artxanda Funicular

Top things to do in Bilbao - Lonely Planet
[www.lonelyplanet.com > Europe > Mediterranean Europe > Spain > Basque Country](http://www.lonelyplanet.com/europe/europe/mediterranean-europe/spain/basque-country) ▾
 The best sights, tours and activities in Bilbao ... Opened in September 1997, Bilbao's shimmering titanium Museo Guggenheim is one of modern architecture's ...

Bilbao Things To Do - Attractions & Must See - VirtualTourist
https://www.virtualtourist.com/travel/.../Bilbao.../Things_To_Do-Bilbao-TG-C-1.html ▾
 Things to Do in Bilbao. Guggenheim Museum & Puppy. Abandoibarra, 2. 73 reviews. Old Town + Cathedral - Casco Viejo & El Arenal. 55 reviews. Zubizuri (White bridge) 10 reviews. Near Bilbao...around Vizcaya / Bizkaia. 15 reviews. Museums. 12 reviews. One hour from Bilbao... Guipuzcoa/Gipuzkoa. 1 review. More Fun things ...

Images for bilbao points of interest 



 3



“Entities” is What a Large Part of Our Knowledge is About

Bilbao photos

Bilbao map

Bilbao
Municipality in Spain

Bilbao, an industrial port city in northern Spain, is surrounded by mountains. It's the de facto capital of Basque Country. It's famous for the Frank Gehry-designed Guggenheim Museum Bilbao, which sparked revitalization when it opened.

Getting there: 1 h 30 min flight, around €200. [View flights](#)

Weather in Bilbao
Weather: 19°C, Wind W at 5 km/h, 79% Humidity

Time in Bilbao
Local time: Monday 9:47 PM

Bilbao Population
Population: 346,574 (2014) Instituto Nacional de Estadística

Colleges and Universities: University of the Basque Country, University of Deusto

Points of interest

Guggenheim Museum	Bilbao Fine Arts	Plaza Nueva	Zubizuri	Doña Casilda

Flights from Paris (all airports) to Bilbao, Spain (BIO)
www.google.fr/flights

Flights from Paris (all airports) to Bilbao, Spain (BIO)

Paris (all airports)	Bilbao, Spain (BIO)
Wed, June 22	Sun, June 26
Non-stop	Air France 1h 35m from €303 Air Europa 1h 35m from €447
Connecting	Multiple airlines 1h 35m+ from €181 Iberia 3h 35m+ from €186 Other airlines 1h 35m+ from €210

[More Google flight results »](#)

Popular schools in Bilbao

Bilbao places to go

4



Push/Pull Techniques for Web Content



Keyword search

Recommendations

User Intent

Augmented search

Domain-aware matching

Semantic search

<https://www.youtube.com/watch?v=L9CP0ltkcJA>

5



Core Entities

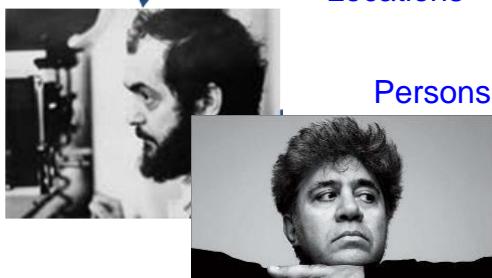
Fall 2019



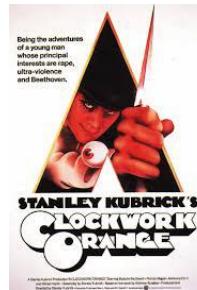
Locations



Organizations



Persons



Movies



6



What Is A Knowledge Base (KB)?

Fall 2019

DBpedia

dbpedia:Stanley_Kubrick	
dbo:birth	dbpedia:Manhattan
hPlace	
rdf:type	foaf:Person
rdf:type	yago:AmericanFilm
	Directors
rdf:type	yago:Amateur
	ChessPlayers

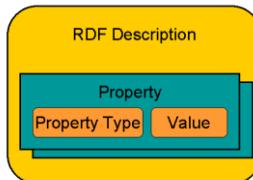
relationships

- Comprehensive, machine-readable descriptions of real-world entities are hosted in knowledge bases (KB)
 - Entity names, types, attributes, relationships, provenance info
- Entities are described as instances of one or several conceptual types and may be linked through relationships
 - The core Semantic Web data model

dbpedia:A_Clockwork_Orange_(film)	
dbo:director	dbpedia:Stanley_Kubrick
dbo:Work/runtime	"136"
foaf:name	"A Clockwork Orange"



facts



7



Knowledge Bases: Scope

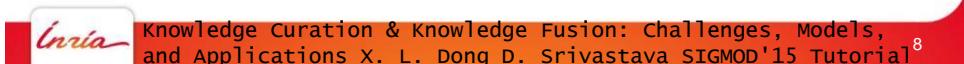
- **Domain-specific Knowledge Base**

- Focus is on a well-defined domain:
 - IMDB for movies, Music-Brainz for music, GeoNames for geo, CIA World Factbook for demographics, etc.



- **Global Knowledge Base**

- Covers a variety of knowledge across domains
 - **Concept-intensive (Intensional)**: Cyc, WordNet
 - **Entity-intensive (Extensional)**: DBpedia, Yago, Freebase, Knowledge Graph, Satori Bing, Knowledge Vault



Knowledge Bases: History

Cyc WordNet



from humans
for humans

$\text{guitarist} \subset \{\text{player}, \text{musician}\}$

$\subset \text{artist}$

algebraist

mathematician

scientist

$\forall x: \text{human}(x) \Rightarrow$

$(\exists y: \text{mother}(x,y)) \wedge$

$\exists z: \text{father}(x,z))$

$\forall x,u,w: (\text{mother}(x,u) \wedge$

$\text{mother}(x,w))$

$\Rightarrow u=w)$

Wikipedia



4.5 Mio. English articles

20 Mio. contributors

WolframAlpha computational knowledge engine

from algorithms
for machines



Knowledge Graph



1985 1990 2000 2005 2010



Adapted from Knowledge Bases in the Age of Big Data Analytics F. Suchanek, Weikum VLDB 2014 Tutorial



Freebase

Freebase is a large knowledge graph containing entities from various domains. It is represented by a central rock labeled "KG" (Knowledge Graph) surrounded by logos of partner websites.

- 46.3M entities
- 1.5K entity types
- 2.67B triples
- 4.5K properties

Logos of partner websites include:

- Baseball Almanac
- clicker
- daylife
- The New York Times
- FANDANGO
- WIKIPEDIA
- TheTVDB.com
- myspace
- NETFLIX
- TV RAGE
- IMDb
- the elderdale project
- facebook
- hulu
- twitter
- NASCAR.COM
- the guardian

Inria John Giannandrea, Freebase – A Rosetta Stone for Entities 11

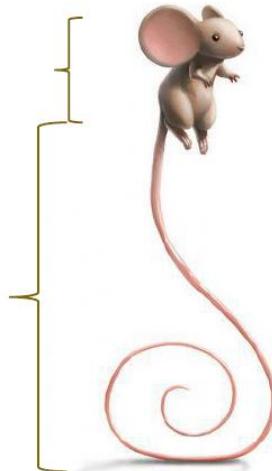


The Long Tail of Knowledge

- Freebase is large, but still very incomplete:

Relation	% unknown in Freebase
Profession	68%
Place of birth	71%
Nationality	75%
Education	91%
Spouse	92%
Parents	94%

- We need automatic KB construction methods



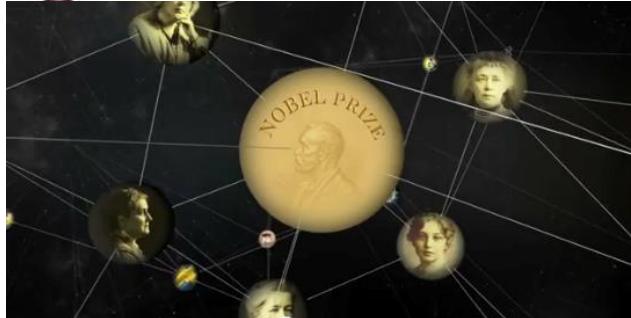
<http://www.flickr.com/photos/sandrei/4691045841/>



Kevin Murphy From Big Data to Big Knowledge CIKM 2013 12



Google's Knowledge Graph (GKG)



- 600M nodes (entities)
- 1.5K entity types
- 20B edges (triples)
- 35K properties

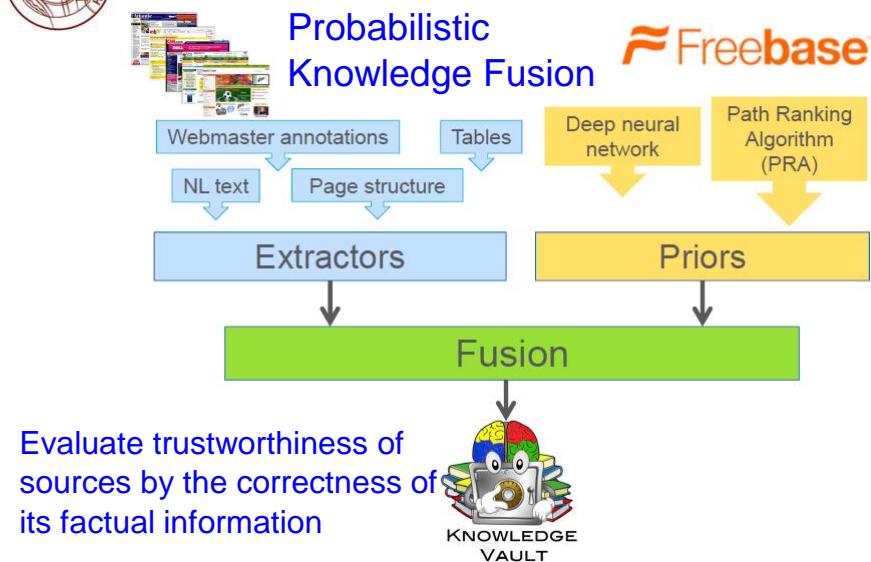
- The Knowledge Graph is used by Google to enhance its search engine's search results with **semantic-search information**
 - GKG is rooted in public KBs such as **Freebase**, **Wikipedia** and the **CIA World Factbook**
 - It's also **augmented at a much larger scale** focusing on a comprehensive breadth and depth of entity descriptions



13



Google Knowledge Vault (GKV)

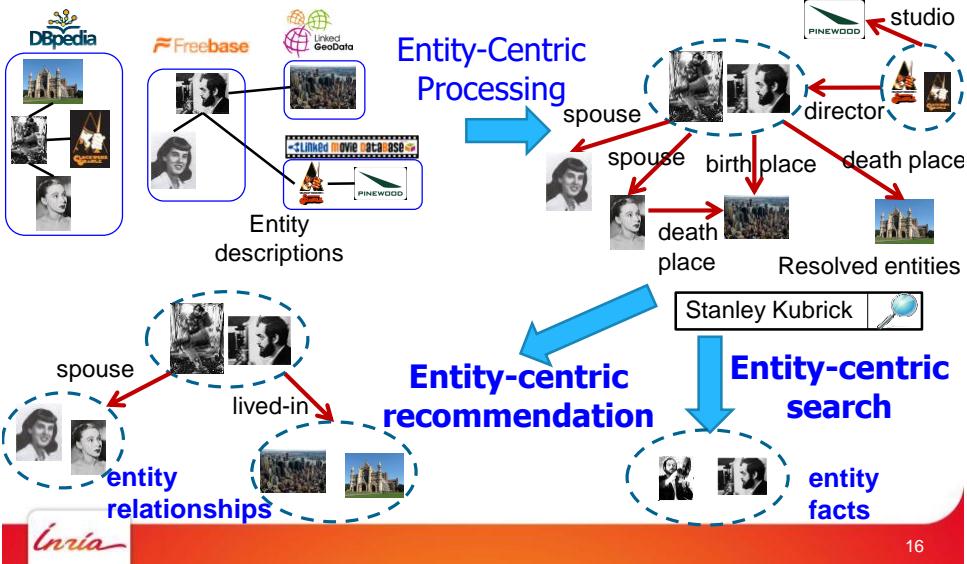




Entity-Centric Applications

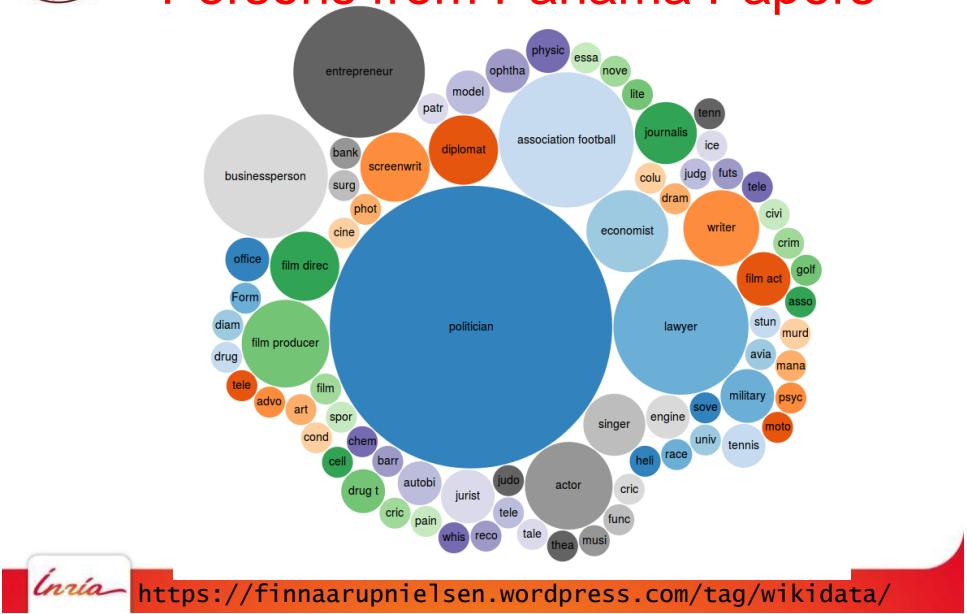
Fall 2019

Combine knowledge regarding an entity from multiple sources to build a rich user experience !



Citizen Journalism: Occupations of Persons from Panama Papers

Fall 2019





Entity-Centric Information Processing

Fall 2019

- Automated construction of entity descriptions
 - *Information extraction*: extract new entities from web/text
 - *Link prediction*: add relationships among entities

- Entity integration and resolution
 - *Knowledge base integration*: instance & ontology mappings
 - *Entity resolution*: merging or splitting similar entities

- Entity-centric access interfaces
 - *Augmented search*: interpret the meaning of queries using entities and compute answers based on a knowledge base
 - *Entity-based matching*: recommend new entities given an entity, a user or a query
 - *Entity-centric summarization*: of textual posts in social media



18



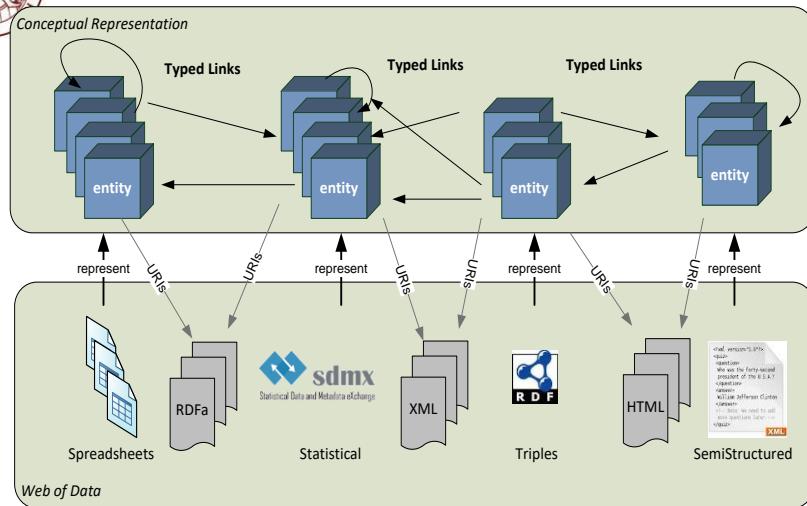
Fall 2019

Describing and Linking Entities: Linked Knowledge Bases & The Web of Data





The Web of Data



A Web of things in the world (aka **entities**),
described by data on the Web

 adapted from chris bizer, richard cyganiak, tom heath ²⁰ <http://linkeddata.org/guides-and-tutorials>



From Open to Linked Data





The Linked Data Principles

Fall 2019

- Anyone can publish data on the Web regarding real-world entities by respecting a minimal set of syntactic conventions
 - Use URIs as names for things
 - Use HTTP URIs so that people machines can look up those names
 - When someone looks up a URI, provide a useful description
 - Include links to other URIs, so that they can discover more things
- Data becomes self-describing
 - Applications encountering data described by an unfamiliar vocabulary, they can resolve its URIs and understand the vocabulary terms by their RDFS or OWL definitions

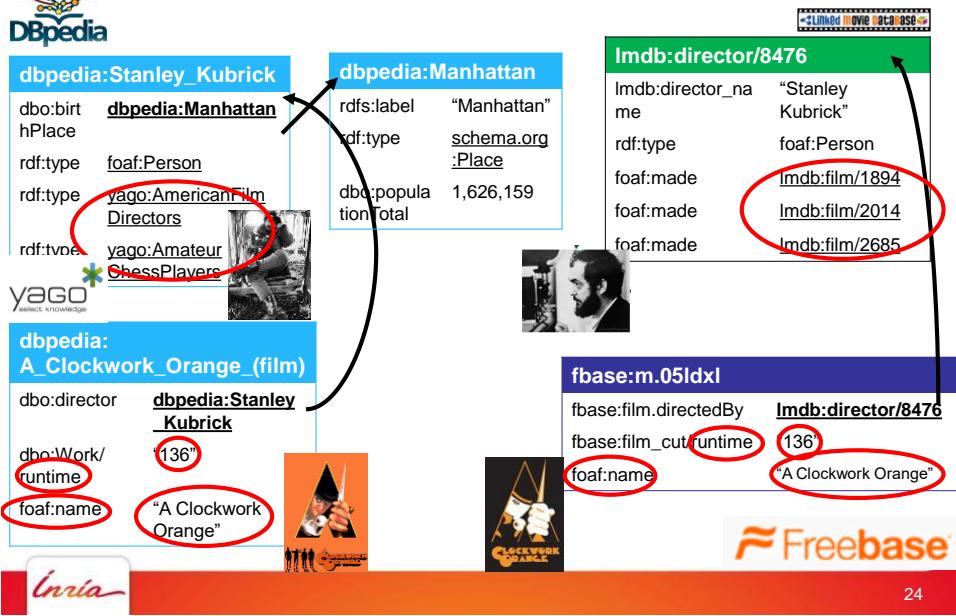


22



Linking Entities in KBs

Fall 2019



24



From Data Silos to the Web of Data

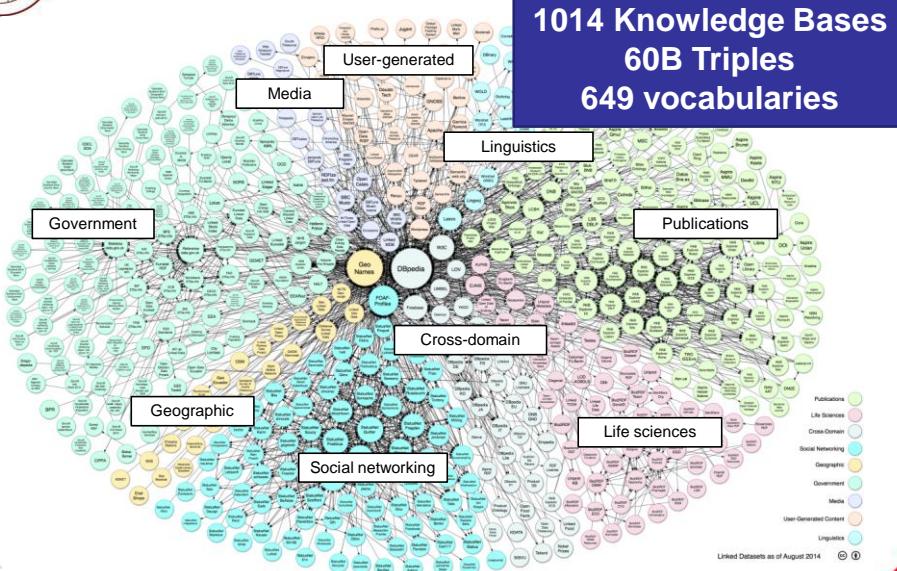


- Thanks to W3C standards we can support:
 - A uniform and universal access to information
 - Easy linking and re-usage of information in different contexts
 - Network effects in adding value to information
- Data Interoperability at the Web Scale
 - Connecting data from diverse sources and across domains
- Openness
 - Many common things are represented in multiple sources
 - Discover new sources at run-time by following links



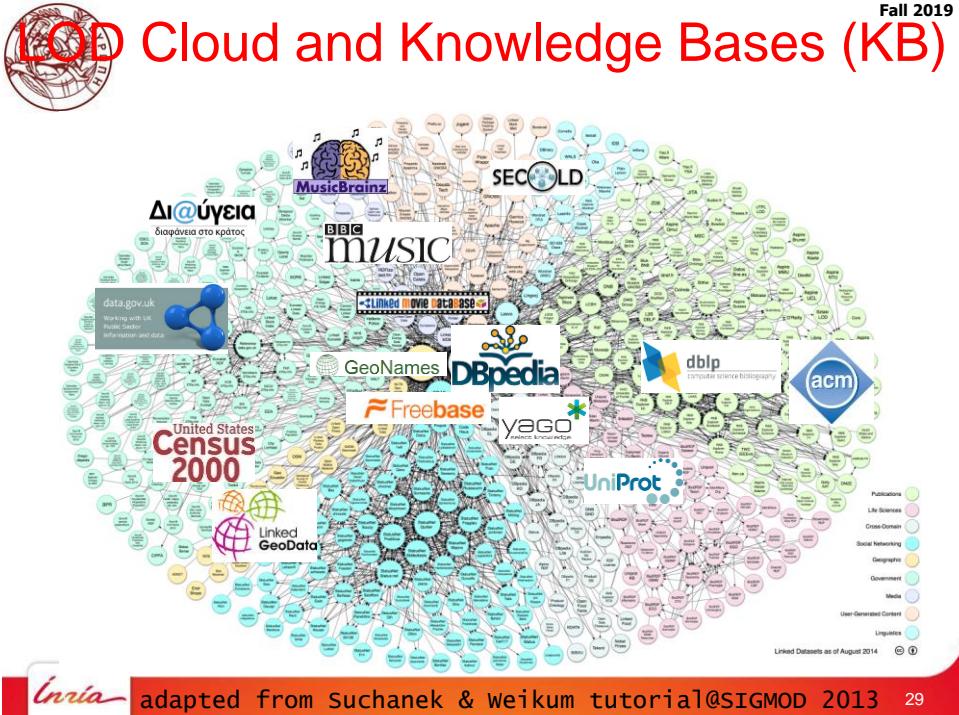
The Link Open Data (LOD) Cloud

1014 Knowledge Bases
60B Triples
649 vocabularies

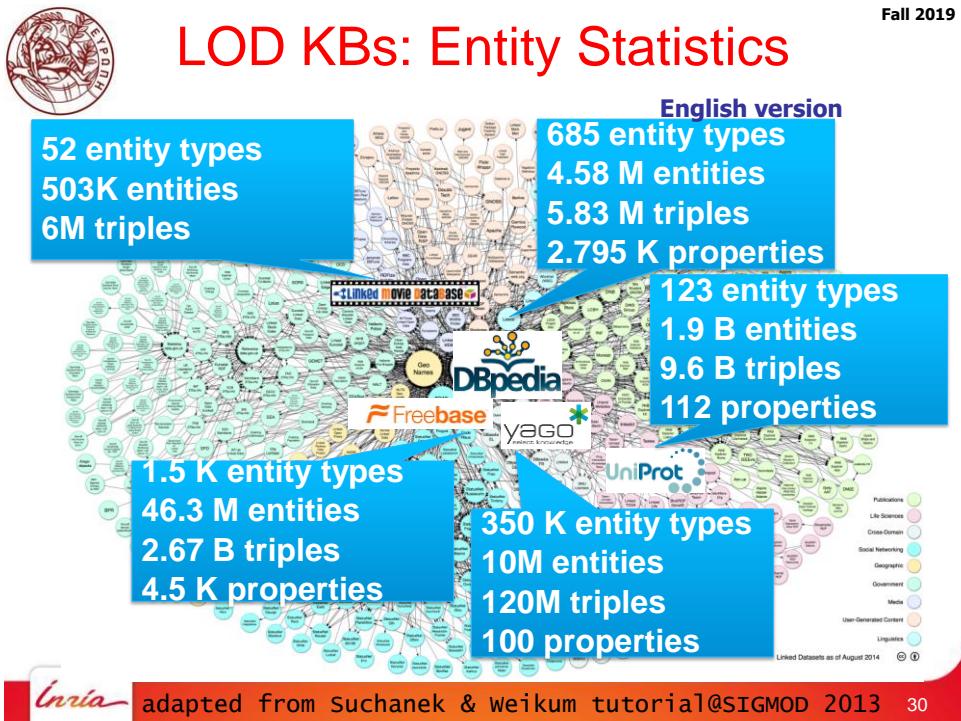


Last updated: 2014-08-30 <http://lod-cloud.net/>

LOD Cloud and Knowledge Bases (KB)



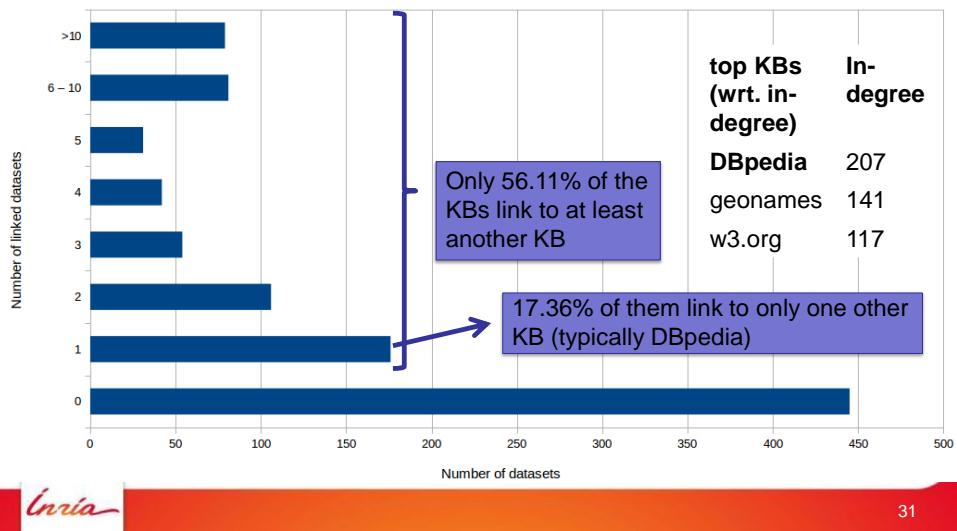
LOD KBs: Entity Statistics





LOD KBs: Entity Interlinking

The LOD cloud diagram is **sparse**



31



Describing and Linking Entities: Description Quality & Entity Resolution





Quality of Entity Descriptions in the Web of Data

Fall 2019

- Given the **open** and **decentralized** nature of the Web, reliability and **usability** of entity descriptions need to be **constantly improved**
 - Incompleteness:** real world entities are only partially described in KBs
 - Redundancy:** descriptions of the same real world entities usually overlap in multiple KBs
 - Inconsistency:** real world entities may have conflicting descriptions across KBs
 - Incorrectness:** errors can be propagated from one KB to the other due to manual copying or automated extraction/fusion techniques



35



Various Forms of KBs Overlapping

Fall 2019

- among KBs (**inter-overlapping**, due to common data sources)



not identical
descriptions,
even if they
stem from the
same source

- within the same KB (**intra-overlapping**, due to data extraction or curation problems)
 - dbpedia:Robert_Soloway and dbpedia:Spam_King refer to the same individual, Robert Soloway, who was nicknamed the “Spam King”
 - less often than inter-overlapping

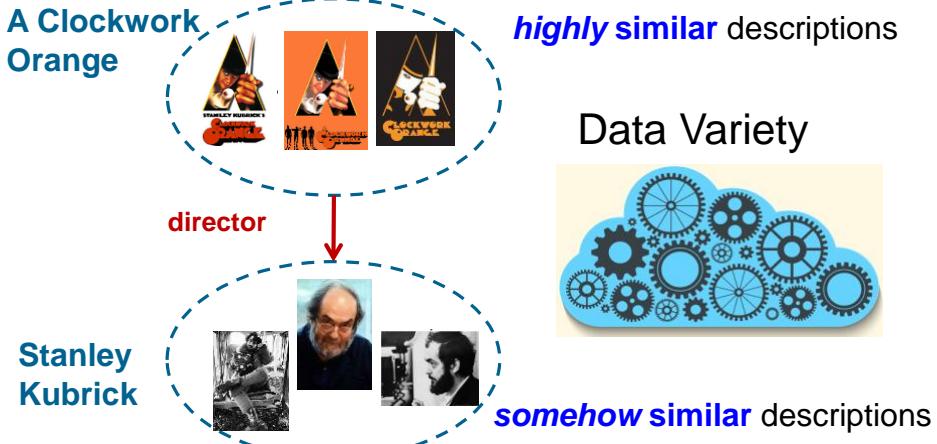


36



Entity Resolution (ER)

The problem of identifying descriptions of the same real-world entity



37



From Highly to Somehow Similar

Highly Similar

- feature many common tokens in the values of semantically related attributes
- heavily interlinked
 - mostly using *owl:sameAs* predicates
- usually provide complementary descriptions of the same entity in terms of facets, temporal evolution
 - good for fusing
- typically met in central KBs

Somehow Similar

- feature significantly fewer common tokens in attributes that are not always semantically related
- sparsely interlinked
 - using various kinds of predicates
- usually complementary descriptions of the same entity in terms of domains, aspects
 - good for linking
- typically met in peripheral KBs



38



Data Variety: Curse or Blessing ?



39



How does ER Improves KB Quality?

- **KB Completeness:**
 - Linking somehow similar descriptions will increase coverage of entity facts and relationships
- **KB Conciseness:**
 - Merging highly similar descriptions will reduce duplicate entity facts and relationships
- **KB Consistency:**
 - Matching similar descriptions will enable to detect conflicting assertions
- **KB Correctness:**
 - Splitting complex descriptions will facilitate entity repairing



40



Web-Scale ER: Challenges



- The two core ER problems, how can we
 - ① effectively compute *entity similarity*
 - ② efficiently resolve sets of entities within or across KBs
- are challenged by the
 - important number of KBs (~ hundreds)
 - large number of entity types & properties (~ thousands)
 - massive volume of entities (~millions)
- Large-scale, progressive, multi-type, cross-domain ER
 - Big Data Volume, Velocity and Variety!



Entity Resolution – Formal Definition

- **Entity resolution:** The problem of identifying descriptions of the same entity within or *across* sources

- $E = \{e_1, \dots, e_m\}$ is a **set** of entity descriptions
- $M : E \times E \rightarrow \{\text{true}, \text{false}\}$ is a **match function**
- The resolution of entities in E results in a partition

$P = \{p_1, \dots, p_n\}$ of E , such that:

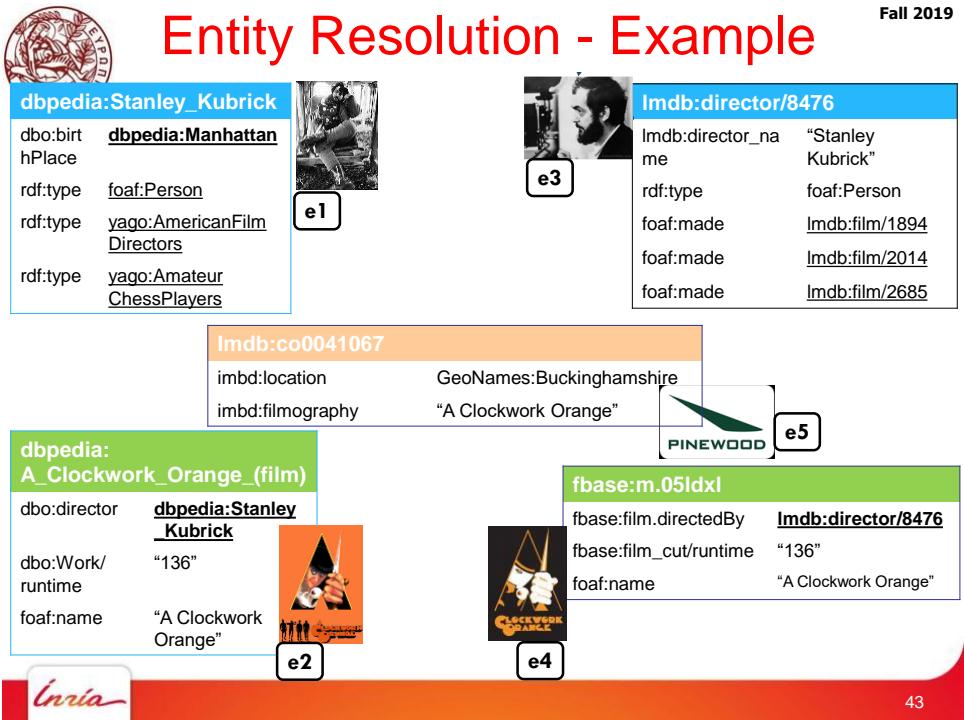
1. $\forall e_i, e_j \in E : M(e_i, e_j) = \text{true}, \exists p_k \in P : e_i, e_j \in p_k$
2. $\forall p_k \in P, \forall e_i, e_j \in p_k, M(e_i, e_j) = \text{true}$

each partition contains only matching descriptions

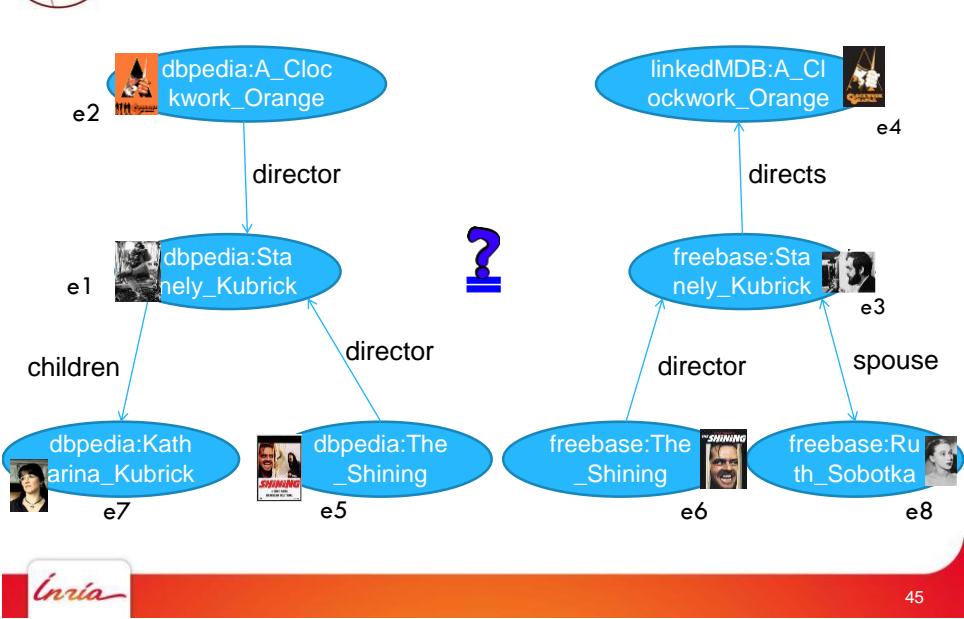
all the matching descriptions are in the same partition



Entity Resolution - Example



Matching Graph-Structured Entities





Single (Pairwise) Entity Matching

Fall 2019

Based on Content Similarity

Matching decisions are **independent**

e1		sim _c (e1,e3) = Jaccard ({Manhattan, Person, AmericanFilmDirectors, AmateurChessPlayers}, {Stanley, Kubrick, Person, 1894, 2014, 2685}) = 0.1	e3
birthPlace	<u>Manhattan</u>		
type	Person		
type	<u>AmericanFilmDirectors</u>		
type	Amateur ChessPlayers		
		thresh = 0.5 	

sim_c: let the content similarity of two descriptions
be the Jaccard similarity of their values' token sets

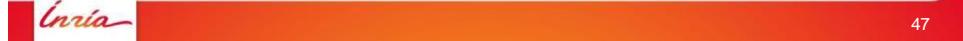
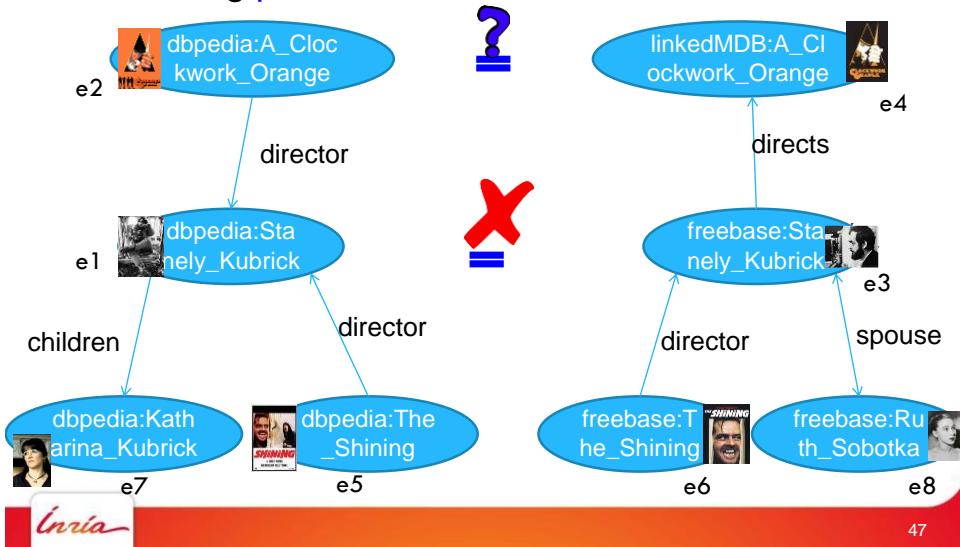


Collective (Joint) Entity Matching

Fall 2019

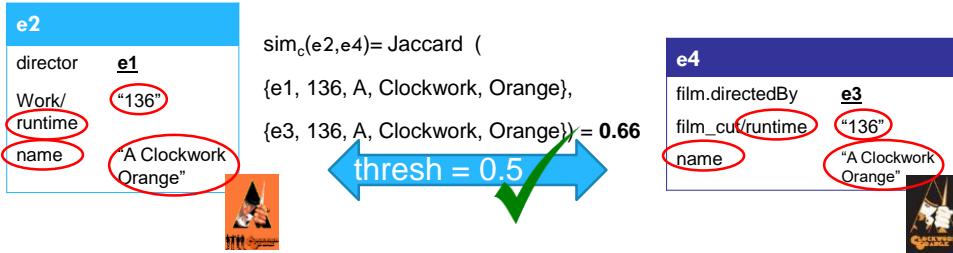
Based on Structural Similarity

One matching provides evidence for another





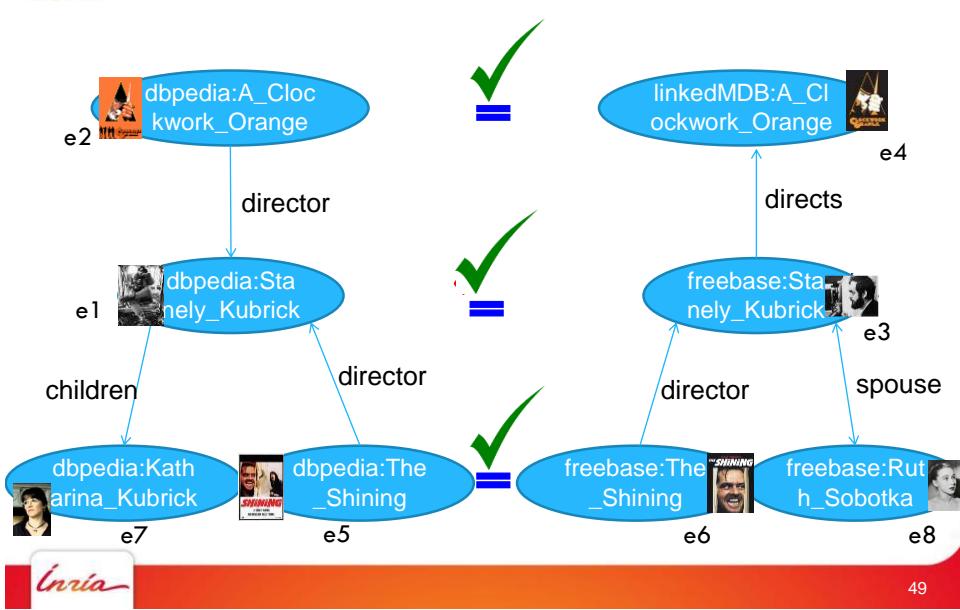
Matching Neighborhood



48



Similarity Propagation

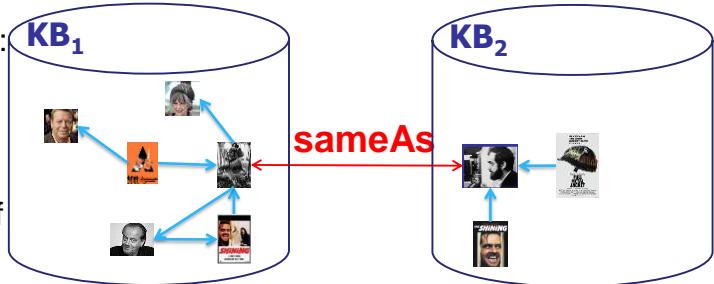


49

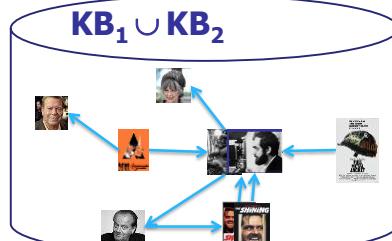


Forms of Entity Resolution (ER)

Record Linkage:
ER without results merging
exploits **exclusivity** of matches



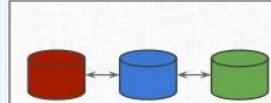
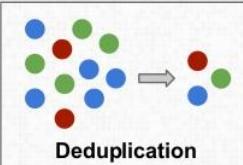
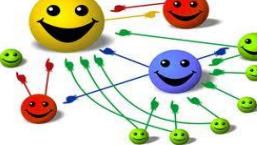
Record Deduplication:
ER with results merging
exploits **transitivity** of matches



50



ER Forms & Similarity

	High similarity in structure	Low similarity in structure
High similarity in content		 Record Linkage
Low similarity in content	 Deduplication	 att & value sim. in a network of relations

string sim. in the values of specific attrs from **one** relation

att & value sim. in a network of relations

set sim. in the values of specific attrs from **two** relations



51



ER Settings

Fall 2019

Two kinds of entity collections as input:

- **Clean**: redundant-free
- **Dirty**: contains redundant entity descriptions

An Entity Resolution (ER) task with input two entity collections can be:

- **Clean-Clean ER**: Given two clean, but overlapping entity collections, identify the common entity descriptions
 - a.k.a. the *record linkage* in databases
- **Dirty-Dirty ER**: Identify unique entity descriptions contained in the union of two dirty input entity collections
 - aka the *deduplication* problem in databases
- **Dirty-Clean Entity Resolution**



52



Comparison with Data Integration Tasks

Fall 2019



F. Naumann, A. Bilke, J. Bleiholder, and M. Weis. Data fusion in three steps: Resolving inconsistencies at schema-, tuple-, and value-level. *IEEE Data Engineering Bulletin*, 29(2):21–31, 2006

53



Scope of the Seminar

Fall 2019

▪ Describing and Linking Entities

- Entity-Centric Applications & Knowledge Bases
- Linked Knowledge Bases & The Web of Data
- Description Quality & Entity Resolution

▪ Matching and Resolving Entities (I)

- Entity Similarity (Content & Context)
- Blocking Techniques (Token, Attribute & URI)
- Block Post-Processing

▪ Matching and Resolving Entities (II)

- Iterative Resolution Techniques (Merging & Matching)
- Progressive Resolution Techniques
- Conclusions & Open Issues



54



Fall 2019

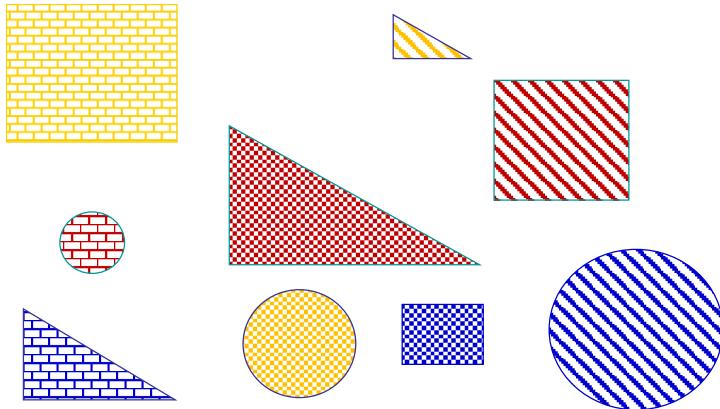
Matching and Resolving Entities (I): Entity Similarity





Similar or Dissimilar?

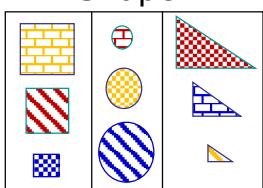
Fall 2019



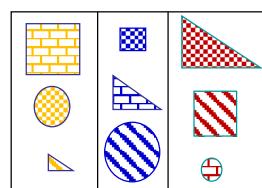
Similarity is Domain Specific

Fall 2019

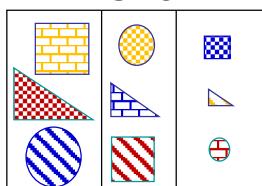
Shape



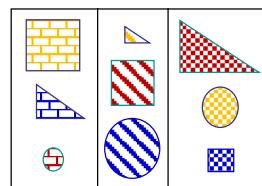
Color



Size



Pattern





Entity Resolution - Match

- **Matches:** Sets of entity descriptions that refer to the same real-world entity, intuitively:
 - Matching descriptions are placed in the same partition
 - All the descriptions of the same partition match
- A **match function** $M()$ maps each pair of entity descriptions (e_i, e_j) to $\{\text{true}, \text{false}\}$
 - $M(e_i, e_j) = \text{true} \Rightarrow e_i, e_j$ are **matches**
 - $M(e_i, e_j) = \text{false} \Rightarrow e_i, e_j$ are **non-matches**



Match Function: Formal Properties

- The match function $M()$ introduce an **equivalence relation** (`owl:sameAs`) among entity descriptions
 - **Reflexivity:** $\forall e_i \in E \text{ if } M(e_i, e_i) = \text{true}$
 - **Symmetry:** $\forall e_i, e_j \in E \text{ } M(e_i, e_j) = M(e_j, e_i)$
 - **Transitivity:** $\forall e_i, e_j, e_k \in E \text{ if } M(e_i, e_j) = \text{true} \text{ and } M(e_j, e_k) = \text{true} \text{ then } M(e_i, e_k) = \text{true}$
- **Additional properties** are useful when linking & repairing entities
 - **Exclusivity** (1-1 Assumption): $\forall e_i, e_k \in E \text{ and } \forall e_j \in E' \text{ if } M(e_i, e_j) = \text{true} \text{ then } M(e_k, e_j) = \text{false}$
 - **Functional Dependency:** $\forall e_i, e_j, e_k, e_l \in E \text{ where } e_i \Rightarrow e_k \text{ and } e_j \Rightarrow e_l \text{ if } M(e_i, e_j) = \text{true} \text{ then } M(e_k, e_l) = \text{true}$





Entity Resolution - Similarity

Fall 2019

- In practice, the match function is defined via a **similarity function** $\text{sim}()$, measuring how similar two entity descriptions are to each other, according to certain comparison criteria
 - Given a **similarity threshold** θ :
 - $M(e_i, e_j) = \text{true}$, if $\text{sim}(e_i, e_j) \geq \theta$
 - $M(e_i, e_j) = \text{false}$, otherwise
 - ML techniques for automatically learning similarity measures are challenged by a Web-scale entity resolution [Köpcke et al. 2010]
 - adaptive learning techniques require training data for each domain [Tejada et al. 2002] [Bilenco et al. 2003]
 - active learning techniques (threshold-based Boolean functions or linear classifiers) work well with highly similar descriptions [Arasu et al. 2010]



Entity Similarity: Example

Fall 2019

dbpedia: A_Clockwork_Orange_(film)	
dbo:director	dbpedia:Stanley Kubrick
dbo:Work/ runtime	136"
foaf:name	"A Clockwork Orange"

	fbase:m.05ldxl	
	fbase:film.directedBy	<u>Imdb:director/8476</u>
	fbase:film_cut/runtime	136
	foaf:name	"A Clockwork Orange"

although not identical e_2 and e_4 are highly similar

dbpedia: <u>Stanley_Kubrick</u>
dbo:birt dbpedia:Manhattan
hPlace dbpedia:New_York_City
rdf:type foaf:Person
rdf:type yago:AmericanFilm
Directors
rdf:type yago:Amateur
ChessPlayers

Imdb:director/8476	
Imdb:director_na	"Stanley Kubrick"
e3	
rdf:type	foaf:Person
foaf:made	Imdb:film/1894
foaf:made	Imdb:film/2014
foaf:made	Imdb:film/2685

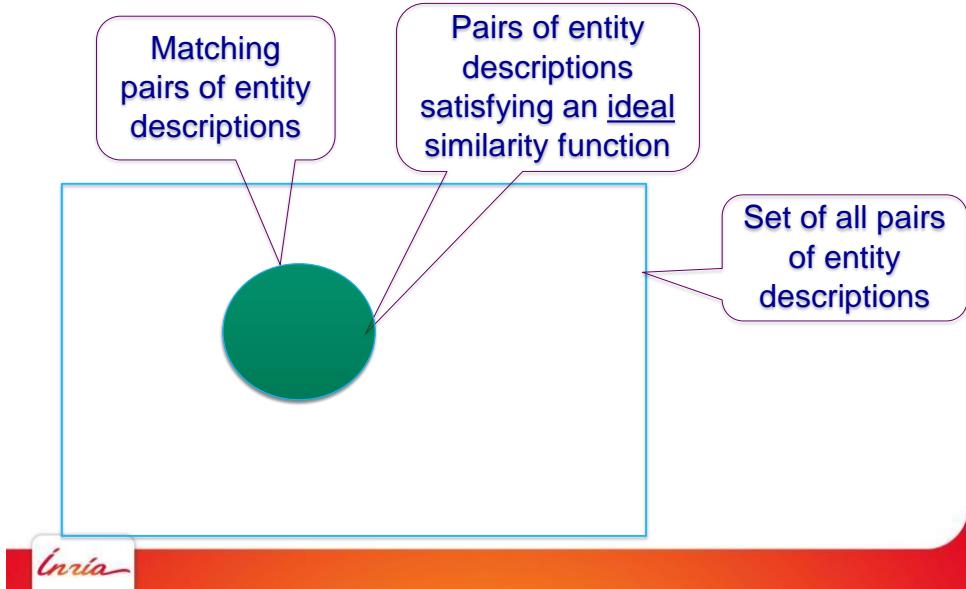
e_1 and e_3 are at best somehow similar





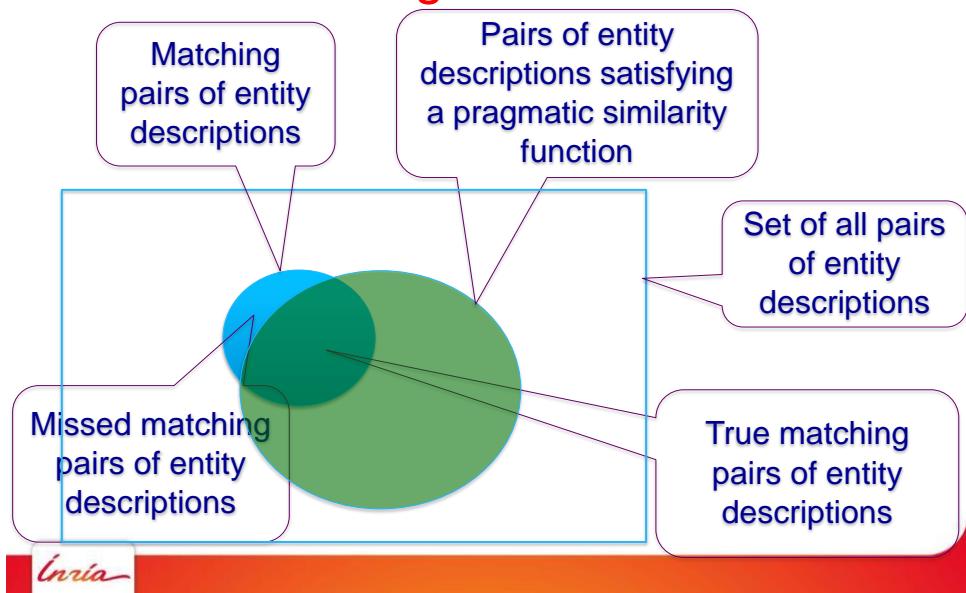
The Role of Similarity Functions: Ideal Case

Fall 2019



The Role of Similarity Functions: A Pragmatic Case

Fall 2019





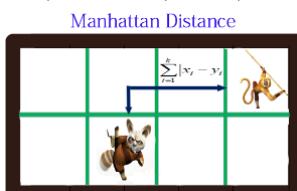
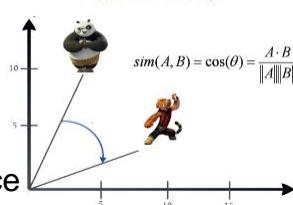
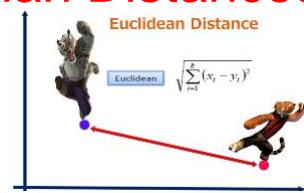
Entity Similarity: Formal Description

- A similarity function $\text{sim}(): E \times E \rightarrow R$ is a **metric**, if for any $e_i, e_j, e_l \in E$ for a given set E , it satisfies the following conditions:
 - $\text{sim}(e_i, e_i) \geq 0$ (**positive**),
 - $\text{sim}(e_i, e_i) \geq \text{sim}(e_i, e_j)$,
 - $\text{sim}(e_i, e_j) = \text{sim}(e_j, e_i)$ (**reflexive**),
 - $\text{sim}(e_i, e_i) = \text{sim}(e_j, e_j) = \text{sim}(e_j, e_i) \Rightarrow e_j = e_i$
 - $\text{sim}(e_i, e_j) + \text{sim}(e_j, e_l) \leq \text{sim}(e_i, e_l) + \text{sim}(e_j, e_j)$ (**triangle inequality**)
- **Normalized** similarity: $\text{sim}(e_i, e_i) \in [0, 1]$
 - $\text{sim}(e_i, e_i) = 1$ for exact match
 - $\text{sim}(e_i, e_j) = 0$ for “completely” different e_i and e_j
 - $0 < \text{sim}(e_i, e_j) < 1$ for some **approximate** similarity
- Any **similarity metric** can be transformed to a **distance metric** and vice versa
 - $\text{sim}(e_i, e_j) = 1 - \text{dist}(e_i, e_j)$



Euclidian vs Non-Euclidian Distances

- A **Euclidean space**: has some number of real-valued dimensions and dense points
 - There is notion of average of two points
 - A **Euclidian distance** is based on the locations of points in such space (L_2 -norm, L_1 -norm)
- A **Non-Euclidian distance** is based on properties of points, but not on their “location” in space
 - Jaccard, Cosine, Edit, Hamming distance



Jaccard Similarity
Set A = {}
Set B = {}
 $|A| = 4$ $|B| = 5$ @dataaspirant.com





Matching and Resolving Entities (I): Content-based Entity Similarity



In Search of Entity Similarity Measures

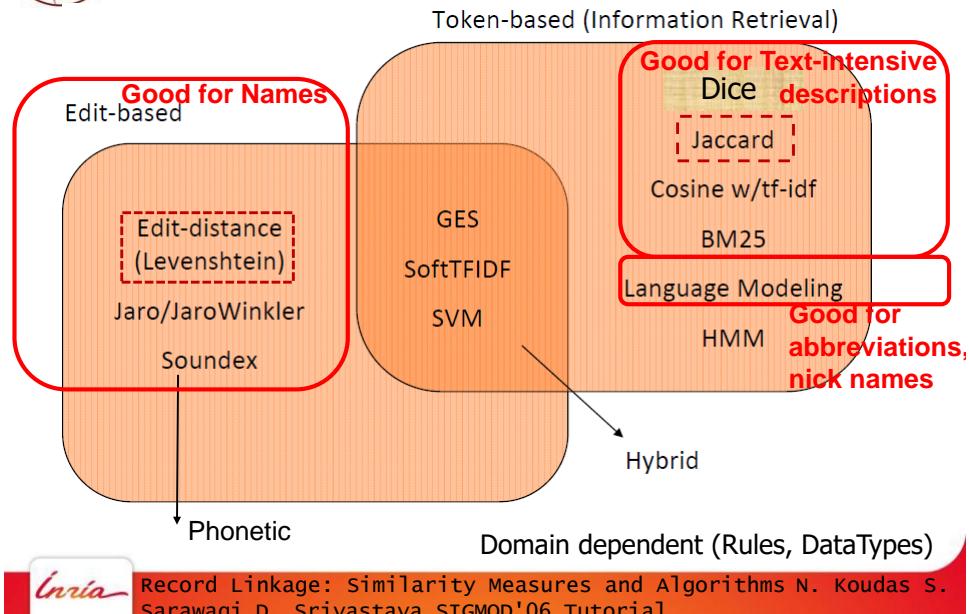
- Defining similarity functions that satisfy the formal properties of metric spaces is, in practice, **too restrictive for non-geometric models**
- **Two main families** of similarity measures for resolving entity descriptions in the Web of data
 - **Content-based**: mostly for measuring string similarity of attribute values in pairs of entity descriptions
 - character-based, token-based
 - **Context-based**: exploit similarity of neighbor descriptions via different entity relationships
 - tree-based, graph-based





String Similarity Measures

Fall 2019



Tokenizing Strings: Separators

Fall 2019

- **Forming words from sequence of characters:**
 - Surprisingly complex in English, can be harder in other languages
- **General idea:** Separate string into tokens using some **separator**
 - Space, hyphen, punctuation, special characters
 - Usually also convert to lower-case
- **Problems:**
 - Sometimes **hyphens** should be considered either as part of the word or a word separator: e.g., spanish-speaking
 - **Apostrophes** can be a part of a word, a part of a possessive, or just a mistake: e.g., master's degree
 - **Periods** can occur in numbers, abbreviations, URLs, ends of sentences, and other situations: e.g., F.M.I





Tokenizing Strings: n-grams

- Split string into short substrings of length n
- Sliding window over string
 - n=2: Bigrams
 - n=3: Trigrams
- Variation: Pad with n – 1 special characters
 - Emphasizes beginning and end of string
- Variation: Include positional information to weight similarities
 - Number of n-grams = $|x| - n + 1$
 - Count how many n-grams are common in both strings

String	Bigrams	Padded bigrams	Positional bigrams	Trigrams
gail	ga, ai, il	⊕g, ga, ai, il, l⊗	(ga,1), (ai,2), (il,3)	gai, ail
gayle	ga, ay, yl, le	⊕g, ga, ay, yl, le, e⊗	(ga,1), (ay,2), (yl,3), (le,4)	gay, ayl, yle
peter	pe, et, te, er	⊕p, pe, et, te, er, r⊗	(pe,1), (et,2), (te,3), (er,4)	pet, etc, ter
pedro	pe, ed, dr, ro	⊕p, pe, ed, dr, ro, o⊗	(pe,1), (ed,2), (dr,3), (ro,4)	ped, edr, dro



Similarity measures Felix Naumann 11.6.2013

75



Token-based Entity Similarity

name	Eiffel Tower
architect	Sauvestre
year	1889
location	Paris e1

about	Eiffel Tower
architect	Sauvestre
year	1889
located	Paris e4

name	Statue of Liberty
architect	Bartholdi Eiffel
year	1886
located	NY e2

about	Lady liberty
architect	Eiffel
location	NY e3
name	White Tower
location	Thessaloniki
year-constructed	1450 e5

$$\text{Jaccard}(\text{tokens}(e_i), \text{tokens}(e_j)) = \frac{|\text{tokens}(e_i) \cap \text{tokens}(e_j)|}{|\text{tokens}(e_i) \cup \text{tokens}(e_j)|}$$

$$\text{Jaccard}(e_1, e_3) = 1/8$$

$$\text{Jaccard}(e_1, e_4) = 1$$

$$\text{Jaccard}(e_1, e_5) = 1/8$$

$$\text{Jaccard}(e_2, e_3) = 3/7$$

$$\text{Jaccard}(e_2, e_4) = 1/11$$

$$\text{Jaccard}(e_2, e_5) = 0/11$$



76

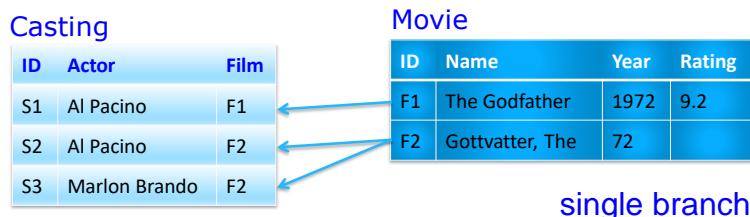


Matching and Resolving Entities (I): Context-based Entity Similarity

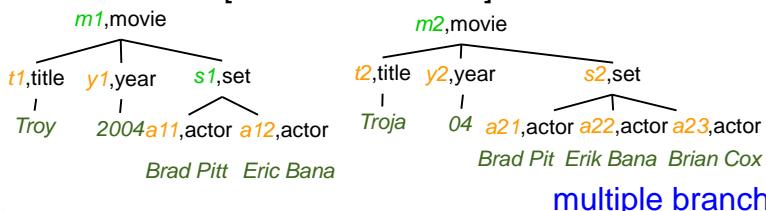


Similarity of Hierarchical Data

- Relational star / snowflake schema [Ananthakrishna et al. 2002]



- Hierarchical XML data [Calado et al. 2010]





DELPHI Containment Metric [Ananthakrishna et al. 2002]

- Hybrid measure considering both similarity of attribute values (*tcm*) and similarity of children sets reached by foreign keys (*fkcm*)
- Similarity of attribute values
 - Divide tuples into tokens → token sets *TS*
 - Compute the edit distance between the tokens of different sets
 - Determine weight of each token via inverse document frequency
 - Thus the idf of a rare term is high, whereas the idf of a frequent term is likely to be low
 - The token similarity metric *tcm* measures which fraction of one tuple *T* is covered by the other tuple *T'* in a relation of interest R_i

$$tcm(T, T') = \frac{\sum idf(TS(T) \cap TS(T'))}{\sum idf(TS(T))} \quad idf_t = \log \frac{N}{df_t}.$$



79



DELPHI Containment Metric [Ananthakrishna et al. 2002]

- Similarity of children sets
 - The children set *CS* of a tuple *T* includes all tuples of a relation R_j referencing *T* from a relation R_i by means of a foreign key
 - Foreign-key containment metric (*fkcm*) measures at what extent the children set of a tuple *T* is covered by the children set of a tuple *T'*

$$fkcm(T, T') = \frac{|CS(T) \cap CS(T')|}{|CS(T)|}$$



80



DELPHI Containment Metric [Ananthakrishna et al. 2002]

Fall 2019

- Combining **tcm** and **fkcm**

- Both **tcm** and **fkcm** are assigned an IDF weight

- Use of a classification function:

$$pos(x) = \begin{cases} 1 & \text{if } x > 0, \\ -1 & \text{otherwise} \end{cases}$$

- Threshold for **tcm**: s1

- Threshold for **fkcm**: s2

- Classification of pairwise comparison between **T** and **T'** using

$$pos(IDF(TS) * pos(tcm(T, T') - s1) + IDF(CS) * pos(fkcm(T, T') - s2))$$

- If final result equals 1, then match, otherwise non-match



81



DELPHI Containment Metric: Example

Fall 2019

ID	Actor	Film
S1	Al Pacino	F1
S2	Al Pacino	F2
S3	Marlon Brando	F2

ID	Name	Year	Rating
F1	The Godfather	1972	9.2
F2	Gottvatter, The	72	

1. Token sets:

$$TS(F1) = \{\text{The, Godfather, 1972, 9.2}\}$$

$$TS(F2) = \{\text{Gottvatter, The, 72}\}$$

2. Attribute value similarities

The = The, Godfather = Gottvatter,
1972 = 72.

3. Weights

For simplification, we assume all tokens have equal weight

4. Token containment metric

$$tcm(F1, F2) = \frac{3}{4}, tcm(F2, F1) = 1$$

5. Children co-occurrence

$$fkcm(F1, F2) = 1, fkcm(F2, F1) = \frac{1}{2}$$

6. Combination of both metrics ($s1 = s2 = 0.5$, weights = 1)

$$pos(pos(3/4 - 0.5) + pos(1 - 0.5)) = 1$$

→ F1 and F2 match



82



Graph-Context Similarity

- Assign a score $s_{x,y}$ to each pair of entities x and y
 - two descriptions are similar if their structural neighborhood is similar (hide attribute values)
- Local similarity Indices
 - Common neighbors (CN)
 - $s_{xy}^{CN} = |\Gamma(x) \cap \Gamma(y)|$ where $\Gamma(x)$ is the set of neighbors
 - Jaccard: Normalized common neighbors (NCN)
 - $\hat{s}_{xy}^{NCN} = \frac{|\Gamma(x) \cap \Gamma(y)|}{|\Gamma(x) \cup \Gamma(y)|}$
 - Adamic-Adar (AA):
 - $s_{xy}^{AA} = \sum_{z \in \Gamma(x) \cap \Gamma(y)} \frac{1}{\log k_z}$ where k_z is the degree of node z and $\frac{1}{\log k_z}$ counts the uniqueness of z, which is inversely proportional to its number of neighbors

83



Content & Context Similarity: Linda [Böhm et al. 2012]

- Works on an entity graph constructed from RDF triples having URIs as subject, predicate and object
 - Literals are stored for each entity e as L(e)
- Matches are identified using a hybrid similarity:
 - String similarity (token-based) of their literal values L(e)
 - Contextual similarity (based on in and out neighbors in the entity graph)
- The context $C(n)$ of e is a set of tuples (p_i, e_i, w_i) , where
 - e_i is a neighboring node of e
 - p_i is the label of the relationship between e and e_i
 - w_i is a numeric weight selected to be higher for less frequent and thus the most discriminative context information

84





Contextual Similarity

Fall 2019

The contextual similarity of nodes n and m is:

$\text{context_sim}(n, m) =$

$$\bullet \sum_{(p_i, z_i, w_i) \in C(n)} \max_{(p_j, z_j, w_j) \in C(m)} w_i \cdot x_{z_i, z_j} \cdot \text{sim}(p_i, p_j), \text{if } |C(n)| \leq |C(m)|$$

$$\bullet \sum_{(p_j, z_j, w_j) \in C(m)} \max_{(p_i, z_i, w_i) \in C(n)} w_j \cdot x_{z_i, z_j} \cdot \text{sim}(p_i, p_j), \text{else}$$

$x_{n,m}$ is 1, if n, m are identified as matches, and 0 else and $\text{sim}(p_i, p_j)$ is the string similarity of the predicates of n, m (edit-distance based)

- It counts the number of common or matching neighbors of two descriptions, which are linked to them in a similar way, i.e., using a relationship with a similar name



85



LINDA Hybrid Similarity

Fall 2019

- The similarity score for descriptions e and e' is:

$$\text{sim}^{\text{LINDA}}(e, e') = \text{content_sim}(e, e') + \beta * \text{context_sim}(c(e), c(e')) - \theta$$

where

β controls the contextual influence

θ is used for re-normalization to values around 0

$$\text{content_sim}_0(n, m) = \frac{|N_n \cap N_m|}{\min(|N_n|, |N_m|) + \ln(|N_n| + |N_m| + 1)}$$

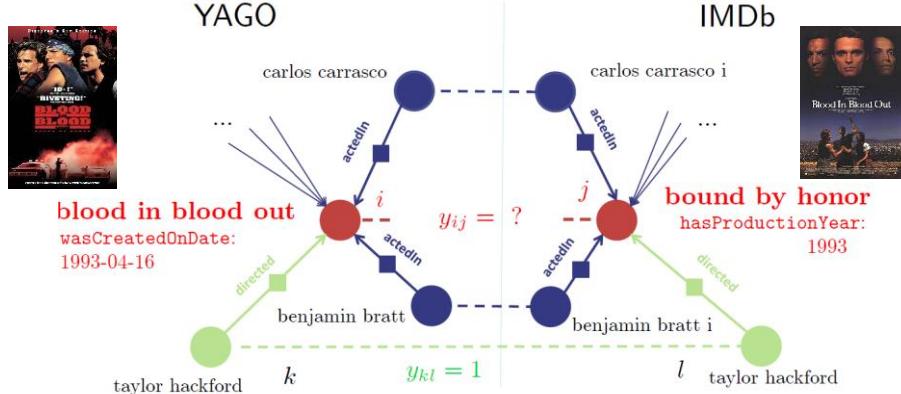
- $\text{sim}^{\text{LINDA}}$ is not a normalized measure as it serves to rank pairs of descriptions based on the evidence that they are matching
 - positive scores reflect likely mappings
 - negative scores imply dissimilarities



86



Content & Structure Similarity: SiGMA [Lacoste-Julien et all 2013]

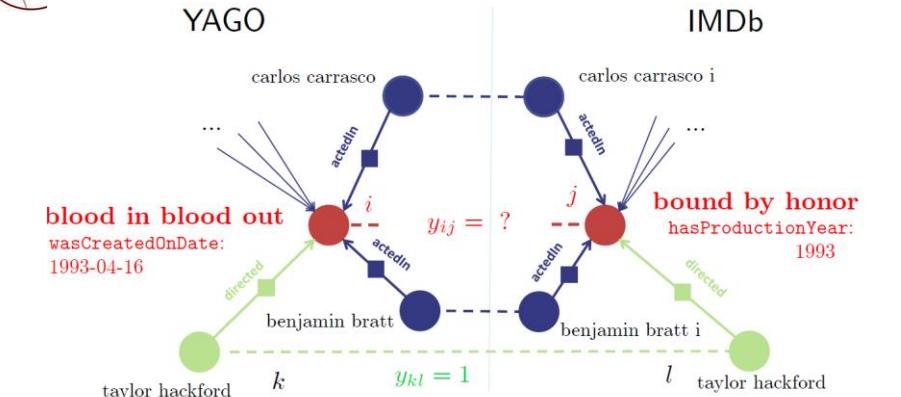


- Even though entities i and j have no tokens in common, the fact that several of their respective neighbors are matched together is a strong evidence that i and j should be matched together
 - Use **neighbors** for scoring and suggesting candidate pairs

87



SiGMA Neighbors Similarity



- Compatible-neighbors N_{ij} : a neighbor k of i being matched to a **compatible** neighbor l of j should encourage i to be matched to j
 - $N_{ij} = \{(k, l) : (i, r, k) \in KB1 \text{ and } (j, s, l) \in KB2 \text{ and relationship } r \text{ is matched to } s\}$

88





SiGMa Similarity Measures

Fall 2019

- **Content similarity:** *static* score of both the string representation of entities (`rdfs:label`) and their other property values

$$s_{ij} = (1 - \beta)\text{string}(i, j) + \beta\text{prop}(i, j) \quad \beta \in [0, 1]$$

- **Context-dependent similarity:** *dynamic* score where the weight $w_{ij,k,l}$ is the contribution of a neighboring matched pair (k, l) to the score of the candidate pair (i, j)

$$\delta g_{ij}(y) \doteq \sum_{(k,l) \in \mathcal{N}_{ij}} y_{kl} (w_{ij,k,l} + w_{kl,i,j})$$

- count the number of compatible neighbors currently matched together for a pair of candidates

$$g_{ij}(y) = \sum_{(k,l) \in \mathcal{N}_{ij}} y_{kl} (\gamma_i w_{ik} + \gamma_j w_{jl})$$

- **Global Score:** $\text{score}(i, j; y) = (1 - \alpha)s_{ij} + \alpha \delta g_{ij}(y)$



89



In Search of Entity Similarity Measures

Fall 2019

- For **highly similar entities** content similarity (i.e., their attribute values) is sufficient
- For **somewhat similar entities** we need to also consider the similarity of the structured context of entities in an **iterative way**
 - Identifying *most discriminating* attributes and relationships is helpful
- An orthogonal issue is the **schematic discrepancy** of attributes and relationships employed in the entity descriptions whose hybrid similarity is assessed
 - **Simple:** either ontological commitments exist or schematic mappings are provided by the users
 - **Complex:** assess similarity of attributes and relationships based on the similarity of their names or values



90

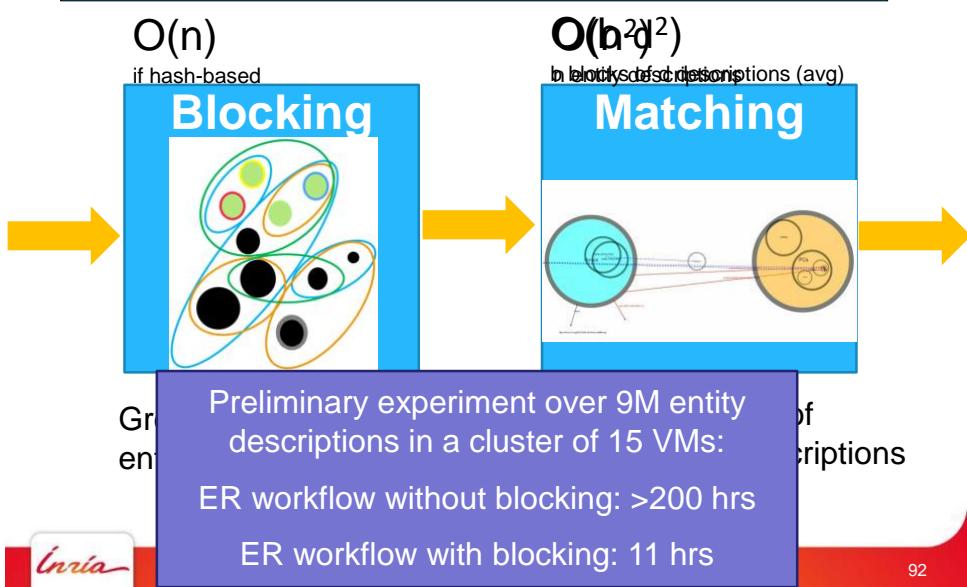


Matching and Resolving Entities (I): Blocking Techniques



Typical ER Workflow

Reduce comparisons leading to non-matching entities





Token Blocking [Papadakis et al. 2011]

- Assume two clean sets KB_1, KB_2 of entity descriptions free of intra-overlapping (*Clean-Clean ER*)
- Each distinct token t_i of values of entity descriptions in $\text{KB}_1 \cup \text{KB}_2$ corresponds to a block
 - Each block contains all entity descriptions sharing the corresponding **token**
 - Pairs originating from the **same (clean) KB** are not compared
- Token blocking offers a **brute-force method** for comparing descriptions even if they are **highly heterogeneous**
 - The same pair of descriptions is contained in many blocks (**redundant comparisons**)
 - Many dissimilar pairs are put in the same block (**unnecessary comparisons**)



Token Blocking Example

name	Eiffel Tower
architect	Sauvestre
year	1889
location	Paris e1

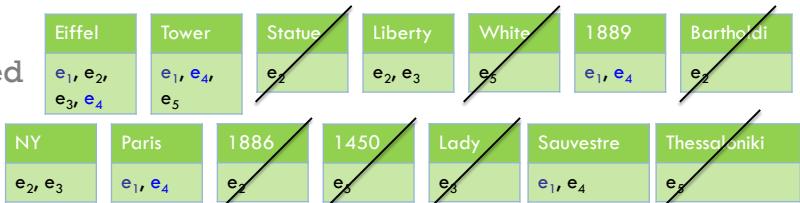
name	Statue of Liberty
architect	Bartholdi Eiffel
year	1886
located	NY e2

about	Lady liberty
architect	Eiffel
location	NY e3

name	White Tower
location	Thessaloniki
year-constructed	1450 e5

Actually, an
**inverted
index**

Generated
Blocks



Redundant comparisons for (e_1, e_4) contained in 5 different blocks!





Token Blocking Example

name	Eiffel Tower
architect	Sauvestre
year	1889
location	Paris e1

name	Statue of Liberty
architect	Bartholdi Eiffel
year	1886
located	NY e2

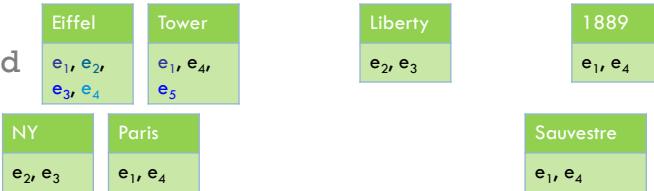
about	Lady liberty
architect	Eiffel
location	NY e3

about	Eiffel Tower
architect	Sauvestre
year	1889
located	Paris e4

Actually, an
inverted
index

name	White Tower
location	Thessaloniki
year-constructed	1450 e5

Generated
Blocks



Unnecessary comparisons between (e₁, e₃), (e₂, e₄), (e₁, e₅)

96



Attribute Clustering Blocking [Papadakis et al. 2013]

- Token blocking totally ignores the semantics of attributes
 - When attribute mappings are not known, attribute clustering considers **similarity of attributes**
 - Attribute similarity is computed w.r.t. the **string similarities of their values**
- Two main steps:
 - Similar attributes are placed together in **non-overlapping clusters**
 - **Token blocking** is performed on the descriptions of each cluster



97

Attribute Clustering Blocking: Example

about	Eiffel Tower	about	Statue of Liberty	about	Auguste Bartholdi	about	Joan Tower
architect	Sauvestre	architect	Bartholdi	born	1834 e13	born	1938 e14
year	1889		Eiffel				
located	Paris e11	year	1886	work	Eiffel Tower	work	Bartholdi Fountain
work	Lady Liberty	located	NY e12	year-constructed	1889	year-constructed	1876
artist	Bartholdi			location	Paris	location	Washington D.C.
location	NY e15				e16		e17

Finding the attribute of KB_2 that is the most similar to “about” of KB_1 :

values {Eiffel, Tower, Statue, Liberty, Auguste, Bartholdi, Joan}

compared to (with Jaccard similarity) :

values of work: {Lady, Liberty, Eiffel, Tower, Bartholdi, Fountain} \rightarrow Jaccard = 4/9

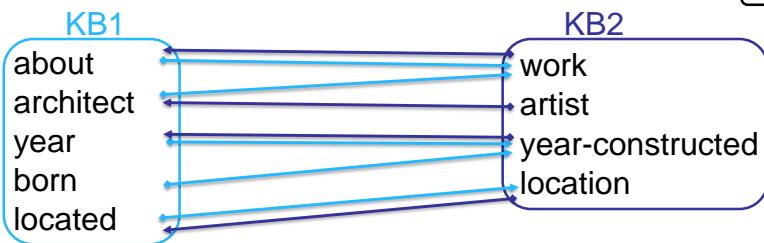
values of artist: {Bartholdi} \rightarrow Jaccard = 1/8

values of location: {NY, Paris, Washington, D.C.} \rightarrow Jaccard = 0

values of year-constructed: {1889, 1876} \rightarrow Jaccard = 0

Attribute Clustering Blocking: Example

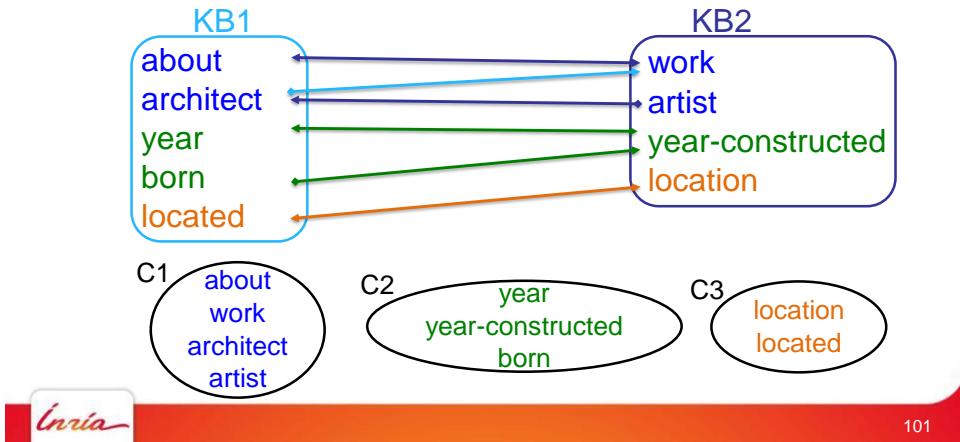
about	Eiffel Tower	about	Statue of Liberty	about	Auguste Bartholdi	about	Joan Tower
architect	Sauvestre	architect	Bartholdi	born	1834 e13	born	1938 e14
year	1889		Eiffel				
located	Paris e11	year	1886	work	Eiffel Tower	work	Bartholdi Fountain
work	Lady Liberty	located	NY e12	year-constructed	1889	year-constructed	1876
artist	Bartholdi			location	Paris	location	Washington D.C.
location	NY e15				e16		e17





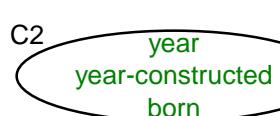
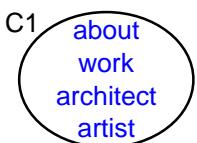
Attribute Clustering Blocking: Example

- Compute the **transitive closure** of the generated attribute pairs
 - Connected attributes form **clusters**
- Example: Pairs (about, work), (work, about), (artist, architect), (architect, work)



Token Blocking in Each Attribute Cluster

about	Eiffel Tower	about	Statue of Liberty	about	Auguste Bartholdi	about	Joan Tower
architect	Sauvestre	architect	Bartholdi Eiffel	born	1834 e13	born	1938 e14
year	1889	year	1886	work	Eiffel Tower	work	Bartholdi Fountain
located	Paris e11	located	NY e12	year-constructed	1889	year-constructed	1876
work	Lady Liberty			location	Paris	location	Washington D.C.
artist	Bartholdi				e16		e17
location	NY e15						



C3.NY	C1.Tower	C1.Bartholdi	→ compare Lady Liberty to Auguste Bartholdi
e ₁₂ , e ₁₅	e ₁₁ , e ₁₄ , e ₁₆	e ₁₂ , e ₁₃ , e ₁₅ , e ₁₇	

Inria

102



Prefix-Infix(-Suffix) [Papadakis et al.12]

- Not only the content but also the **identity (URI)** of descriptions reveal valuable semantics regarding similarity of entities
 - E.g., 66% of the 182 million URIs of BTC09 follow the scheme: Prefix-Infix(-Suffix)
- Example: http://en.wikipedia.org/wiki/Linked_data#Principles.html
 - **Prefix** describes the source, i.e. domain, of the URI
 - **Infix** is a local identifier or a named anchor
 - The optional **Suffix** contains details about the format, e.g. .html, .rdf, .nt, etc.
- **Infix blocking**
 - The blocking key is the URI infix of the entity description
- **Infix profile blocking**
 - The blocking keys are the URI infixes of attribute values in entity descriptions



Infix Blocking

yago:Statue_of_Liberty		dbpedia:Statue_of_Liberty		fb:m.072p8		geonames:5139572	
skos:prefLabel	Statue of Liberty	rdfs:label	Statue of Liberty	fb:official_name	Statue of Liberty	geoname:s:name	Statue of Liberty
yago:isLocatedIn	yago:Liberty_Island	dbprop:location	dbpedia:Liberty_Island	fb:contains_d_by	fb:m.026kp2	geoname:s:nearby	geonames:5124330
e1		e2		e3		e4	
yago:Tina_Brown							
skos:prefLabel	Tina Brown	e5		e3		e4	
yago:linksTo		yago:Liberty_Island					

Statue_of_Liberty	m.072p8	5139572	Tina_Brown
e ₁ , e ₂	e ₃	e ₄	e ₅





Infix Profile Blocking

yago:Statue_of_Liberty	dbpedia:Statue_of_Liberty	fb:m.072p8	geonames:5139572
skos:prefLabel	Statue of Liberty	rdfs:label	Statue of Liberty
yago:isLocatedIn	yago:Liberty_Island	dbprop:location	dbpedia:Liberty_Island
			e2
yago:Tina_Brown			
skos:prefLabel	Tina Brown		e3
yago:linksTo	yago:Liberty_Island		

(e₁, e₃) true positive

(e₁, e₅) false positive

Liberty_Island	m.026kp2	5124330
e ₁ , e ₂ , e ₃ , e ₅	e ₃	e ₄

- The effectiveness of these approaches relies on the good naming practices of the KBs publishing entity descriptions



Other Blocking Techniques

- Similarity Join algorithms:** relies on an inverted index from the most discriminating (non-frequent) tokens of the descriptions and the prefix filtering principle [Chaudhuri et al., 2006]

$$\text{Jaccard}(x, y) \geq t \Rightarrow |x \cap y| \geq \frac{t}{1+t} \cdot (|x| + |y|)$$
 - prefix filtering is effective only when the similarity threshold is extremely high
- Frequent itemsets blocking:** build blocks for sets of tokens that frequently co-occur in entity descriptions
 - may significantly reduce the number of candidate pairs
 - may also significantly increase missed matches, especially between descriptions with few common tokens
- Multidimensional blocking:** constructs a collection of blocks for each similarity function used to resolve entities and aggregates them into a single multidimensional collection, by taking into account the similarities of descriptions that share blocks





Placing Entities in the Same Block

Method	Criterion
Token Blocking [Papadakis et al., 2011]	The descriptions have a common token in their values
Attribute Clustering Blocking [Papadakis et al., 2013]	The descriptions have a common token in the values of attributes that have similar values in overall
Prefix-Infix(-Suffix) [Papadakis et al., 2012]	The descriptions have a common token in their literal values, or a common URI infix
ppjoin+ [Vernica et al., 2010]	The descriptions have a common token in their p first tokens (sorted in ascending frequency order)
Frequent itemsets [Kenig and Gal, 2013]	The descriptions have frequently co-occurring tokens in their values



107



Blocking Algorithms Comparison

Method (suffix)	Partiti oning	Overlapping			Algorithm	
		positi ve	negati ve	neut ral	Hash- based	Sort- Based
<i>Canopy Clustering</i> [McCallum et al. 00]			•			•
<i>Standard blocking</i> [Fellegi& Sunter 69, Kolb et al.,12b]	•				•	
<i>Q-grams</i> [Gravano et al., 01]		•				•
<i>Suffixes</i> [Aizawa & Oyama, 05]		•			•	
<i>Sorted neighborhood</i> [Hernández & Stolfo 95, Kolb et al.,12a]				•		•
<i>Adaptive sorted neighborhood</i> [Yan et al., 07]	•					•
<i>Token blocking</i> [Papadakis et al.,11]		•			•	
<i>Attribute clustering</i> [Papadakis et al.,13]			•			•
<i>Prefix-infix (suffix)</i> [Papadakis et al. 12]		•			Schema Agnostic	
<i>ppjoin+</i> [Vernica et al., 10]		•			•	
<i>Frequent itemsets</i> [Kenig&Gal,13]		•			•	



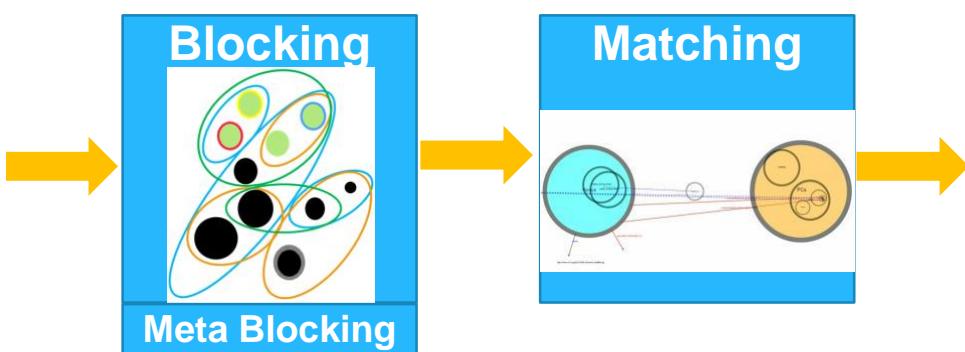


Matching and Resolving Entities (I): Block Post-Processing



Meta-blocking: Improving the Efficiency of Blocking

Eliminate redundant comparisons &
reduce unnecessary comparisons





Meta-blocking – The Blocking Graph

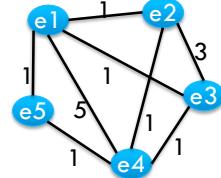
Blocks (ToB):

Eiffel	Tower	Liberty
e ₁ , e ₄ , e ₂ , e ₃	e ₁ , e ₄ , e ₅	e ₂ , e ₃
NY	1889	Paris
e ₂ , e ₃	e ₁ , e ₄	e ₁ , e ₄
Sauvestre		
e ₁ , e ₄		

14 comparisons to identify 2 matches
e₁-e₄ and e₂-e₃

Blocking graph (Nodes: entity descriptions, Edges: common block)

Pruned blocking graph (discard edges with weight below avg.: 1.75)



edge weights = #common blocks

2 comparisons to identify 2 matches

Prune edges to discard unnecessary comparisons between non-matches based on positive overlapping evidence



Edge Weighting and Pruning

- **Weighting Schemes** (how to weight the edges)
 - **Common Blocks (CBS)**: $w_{i,j} = |B_{i,j}|$
 - **Jaccard (JS)**: $w_{i,j} = |B_{i,j}| / (|B_i| + |B_j| - |B_{i,j}|)$
 - **Enhanced CBS (ECBS)**: $w_{i,j} = CBS \cdot \log(|B|/|B_i|) \cdot \log (|B|/|B_j|)$
- **Pruning Methods** (which edges to prune)
 - **WEP**: Keep edges with weight above average
 - **CEP**: Keep top-K edges overall
 - **WNP**: Keep, for each node, the edges with weight above a local average
 - **CNP**: Keep, for each node, its top-K edges

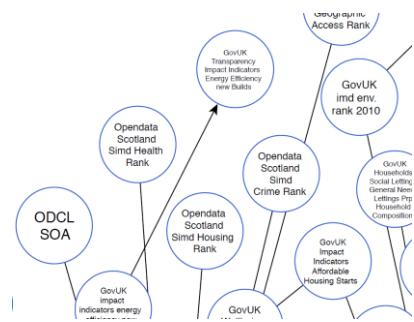
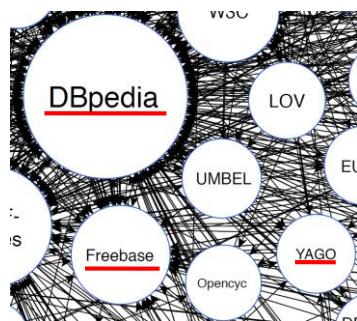




Matching and Resolving Entities (II): Iterative Resolution Techniques



Central vs. Peripheral KBs



In Search of Similarity Evidence in KBs

[Efthymiou et al. 2016, 2015]

- Attribute-based comparisons
 - unique attributes (e.g., rdfs:label) provide strong evidence
 - >90% of matching pairs have >80% overlap similarity in the values of rdfs:label
- Content-based comparisons
 - central KBs: 3-4 common tokens in entity values
 - peripheral KBs: 1-2 common tokens in entity values
 - blocking algorithms miss up to 30% matches in peripheral KBs
- Relationships-based comparisons
 - matching neighbors provide positive evidence
 - >92% of pairs with at least one matching neighbor, are matches in most KBs
 - some types of relationships provide strong negative evidence
 - dissimilar values for wasBornIn indicate a non-matching pair



117

Types of Missed Matches (FNs)

- Type A: a third, matching description (transitivity)

applicable to identify matches within a KB
- Type B: matches of their neighbors

can identify matches both within a KB and across different KBs

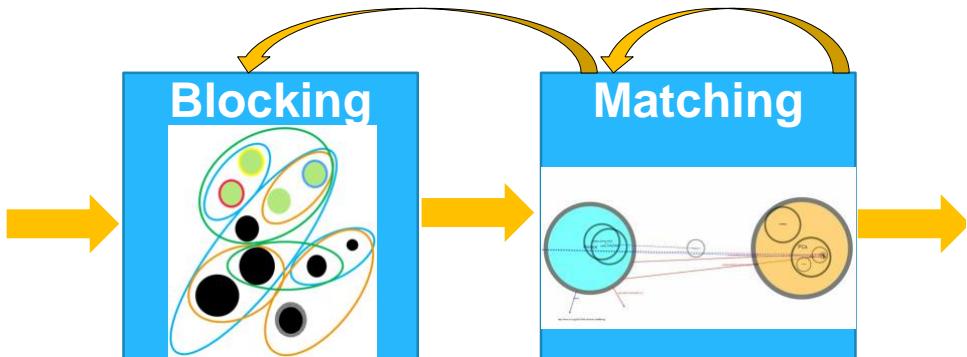


118



Iterative ER

Increase the number of matching entities



Iterative ER: identify new matches based on partial results either of matches or of merges

- Good for high Variety



Iterative ER Approaches

- **Merging-based:** new matches can be found by exploiting merged (more complete) descriptions of previously identified matches
 - Idea: ER resembles a **database self-join operation** (of the initial set of descriptions with itself)
 - But: No knowledge about which descriptions may match, so all pairs of descriptions need to be compared
- **Matching-based:** If descriptions related to entity e_i are matching to descriptions related to e_j , then e_i and e_j are likely to match
 - Idea: ER resembles to a **graph traversal problem** in which **similarity is propagated** until a fixed point is reached
 - Use **positive or negative evidence** for prioritize similarity re-computation





Matching and Resolving Entities (II): Merging-based Iterative Resolution



Merging-based ER– Formal Definition

- Let $E = \{e_1, \dots, e_m\}$ a set of entity descriptions and
 - $M : E \times E \rightarrow \{\text{true, false}\}$ is a match function
 - $\mu : E \times E \rightarrow E$ is a partial merge function
- The merging-based (generic) resolution of entities in E is a set of descriptions E' , such that:
 1. $\forall e_i, e_j \in E : M(e_i, e_j) = \text{true}, \exists e_k \in E' : \mu(e_i, e_j) \leq e_k$
 2. $\forall e_k \in E', \forall e_l \in E, e_l \leq e_k$
 3. no strict subset of E' satisfies conditions 1 and 2

where $e_l \leq e_k$ means that e_k is at least as informative than e_l , regarding the same real-world entity
- Note: $E' \subseteq \hat{E}$ (the merge closure)
 - Condition 1: E' cannot produce more than \hat{E}
 - Condition 2: E' produce at least all information of \hat{E}
 - Condition 3: Minimal solution





Match and Merge Functions: ICAR Properties [Benjelloun et al., 2009]

Fall 2019

- **Idempotence:** a description always matches itself, and merging it with itself still yields the **same description**
 - $\forall e_i \in E$ if $M(e_i, e_i) = \text{true}$ then $\mu(e_i, e_i) = e_i$
- **Commutativity:** **direction** of match and merge is **irrelevant**
 - $\forall e_i, e_j \in E$ if $M(e_i, e_j) = M(e_j, e_i) = \text{true}$ then $\mu(e_i, e_j) = \mu(e_j, e_i)$
- **Associativity:** **order** of merge is **irrelevant**
 - $\forall e_i, e_j, e_k \in E$ if $\mu(\mu(e_i, e_j), e_k)$ and $\mu(e_i, \mu(e_j, e_k))$ exist then
 $\mu(\mu(e_i, e_j), e_k) = \mu(e_i, \mu(e_j, e_k))$
- **Representativity:** merging does **not lose matches**; no “negative evidence”
 - if $e_k = \mu(e_i, e_j)$ then $\forall e_l \in E$ such that $M(e_i, e_l) = \text{true}$ also $M(e_k, e_l) = \text{true}$
- Transitivity is **not assumed!**



123



Merge Domination & Monotonicity

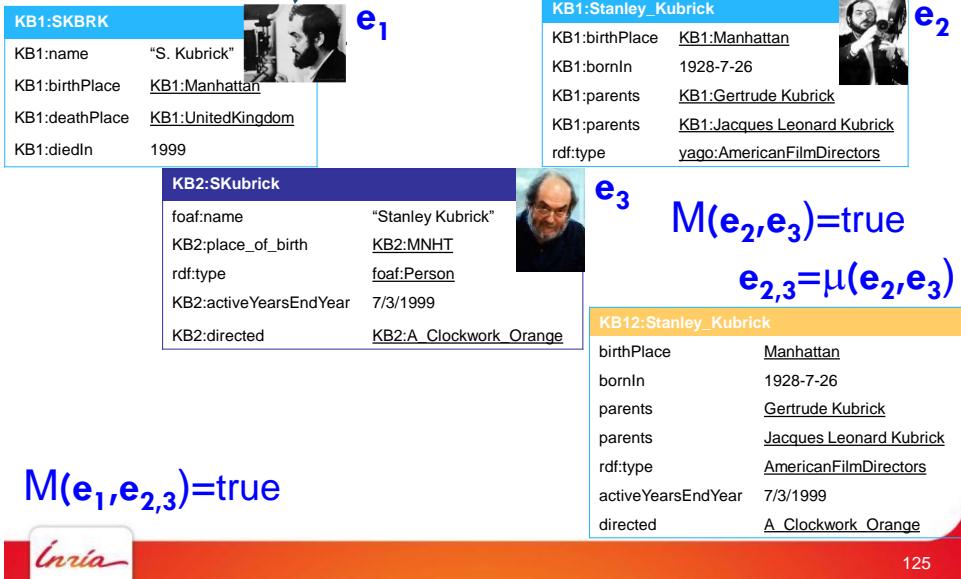
Fall 2019

- When the **match** and **merge** functions satisfy the ICAR properties, there is a **natural domination (partial) order** of descriptions
- Given two descriptions, e_1 and e_2 , we say that e_1 is **merge dominated** by e_2 , denoted $e_1 \leq e_2$, if $M(e_1, e_2) = \text{true}$ and $\mu(e_1, e_2) = e_2$
 - e_1 does not add information
- **Merged descriptions always dominates** the ones they have derived from
 - $\forall e_1, e_2 \in E$ such that $M(e_1, e_2) = \text{true}$, it holds that $e_1 \leq \mu(e_1, e_2)$ and $e_2 \leq \mu(e_1, e_2)$
- **Match function is monotonic**
 - $\forall e_1, e_2, e_3 \in E$ If $e_1 \leq e_2$ and $M(e_1, e_3) = \text{true}$, then $M(e_2, e_3) = \text{true}$
- **Merge function is monotonic**
 - $\forall e_1, e_2, e_3 \in E$ If $e_1 \leq e_2$ and $M(e_1, e_3) = \text{true}$, then $\mu(e_1, e_3) \leq \mu(e_2, e_3)$
- $\forall e_1, e_2, e_3 \in E$ If $e_1 \leq e_3$, $e_2 \leq e_3$ and $M(e_1, e_2) = \text{true}$, then $\mu(e_1, e_2) \leq e_3$



Generic Entity Resolution with Swoosh F. Naumann 20.6.2013 124

What is Best Sequence of Match, Merge Calls that Give us Right Answer?

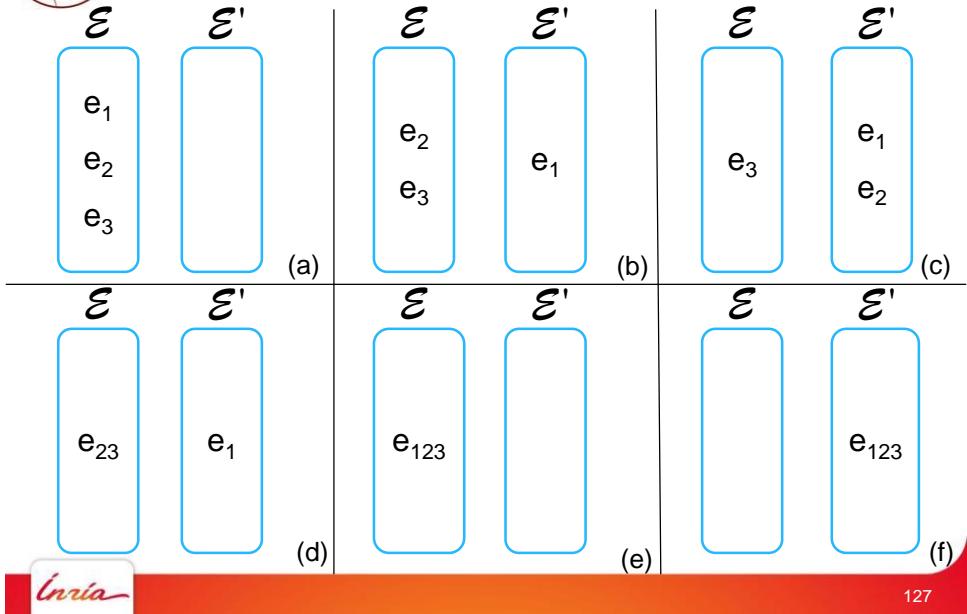


R-Swoosh [Benjelloun et al., 2009]

- Assumes ICAR and merge domination
- Idea 1:** if $M(e_1, e_2) = \text{true}$ we can remove e_1 and e_2
 - Whatever would match e_1 or e_2 now also matches $\mu(e_1, e_2)$
 - Representativity and associativity
- Idea 2:** Removal of dominated descriptions is not necessary as a last step in the algorithm
 - Assume e_1 and e_2 appear in final answer where $e_1 \leq e_2$
 - Then $M(e_1, e_2) = \text{true}$ and $\mu(e_1, e_2) = e_2$
 - Thus comparison of e_1 and e_2 should have generated merged description e_2 , and e_1 should have been eliminated

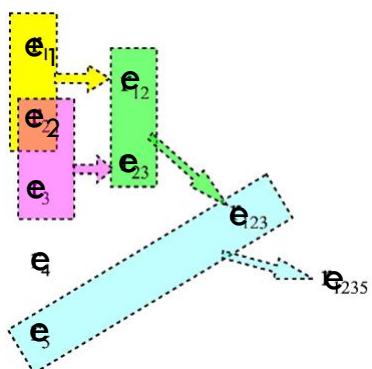


R-Swoosh Example



ER with ICAR properties

- General ER process is guaranteed to be **finite**
- Entity descriptions can be **matched** and **merged** in any order
- Dominated** entity descriptions can be **discarded** anytime
- Union class of match and merge
 - Union-merge:** All values are kept in merged descriptions
 - Union-match:** At least one of the values in common
- Commutativity of **match** and **merge** functions for **highly heterogeneous** descriptions does **not always hold**





Matching and Resolving Entities (II): Matching-based Iterative Resolution



Positive and Negative Match Evidence

Fall 2019

		Evaluated on individual entities	Evaluated on interrelated entities
		Hard Constraints	Soft Constraints
Positive Evidence	FD: if $M(e_i, e_j) = \text{true}$ then $M(e_k, e_l) = \text{true}$	FD: if $M(e_i, e_j) = \text{true}$ then most likely $M(e_k, e_l) = \text{true}$	
	Example: if two movies match then their director also match	Example: if two movies match then their actors are most likely to match	
Negative Evidence	FD: if $M(e_i, e_j) = \text{false}$ then $M(e_k, e_l) = \text{false}$	FD: if $M(e_i, e_j) = \text{false}$ then most likely $M(e_k, e_l) = \text{false}$	
	Example: if two directors don't match then movies directed by them don't also match	Example: if two movies don't match then their actors are less likely to match	

Negative constraints are usually stated by domain experts

Constrains can be **recursive**

Inria

L. Getoor, A. Machanavajjhala Entity Resolution for Big Data
KDD'13 Tutorial

Additional Constraints [Shen et al, 2005]

Type	Example
Aggregate	No director has produced more than four movies per year
Subsumption	If a movie X from Yago matches a movie Y from IMBD then each studio cited by in Y matches some studio cited by in X
Neighborhood	If actors X and Y share similar names and some co-actors in a movie, they are likely to match
Incompatible	No director has filmed a movie in studios both in Africa and Japan
Layout	If two movies with similar names are mentioned by the same review they are likely to match
Key/Uniqueness	Actors of the same movie must refer to distinct persons
Ordering	If two paper citations match, then their authors will be matched in order
Individual	'Victorina Mérida Rojas' matches with 'Victoria Abril'



Matching-based Iterative ER

- **Pair-wise ER:** matching decisions are made independently
 - Deduplication, Linkage
- **Collective ER:** matching decisions depend on other matching decisions according to positive and negative evidence (**general constraints**)
 - **Similarity propagation** approaches (more scalable)
 - Dependency graphs of matching decisions
 - Collective relational clustering
 - **Probabilistic** approaches (scalability is an open issue)
 - Generative Models (for acyclic dependencies between match decisions)
 - Undirected Models based on Markov Networks (sometimes with a first-order logic syntax)
 - **Hybrids** of constraints and probabilistic models





Similarity Propagation Approaches

- A graph structure for encoding the similarity between entity descriptions and matching decisions, and iteratively assess matching of entities by propagating similarity values
 - Details of how the graph is constructed and traversed and how (content and context) similarity is computed vary
- Similarity-propagation ER: the match function is re-computed at each iteration step by considering previous matching decisions:
 - $M^n(e_i, e_j) = \text{true}$, if $\text{sim}^{n-1}(e_i, e_j) \geq \theta$
 - $M^n(e_i, e_j) = \text{false}$, if $\text{sim}^{n-1}(e_i, e_j) \leq \theta'$
 - $M^n(e_i, e_j) = \text{undecided}$, otherwise
- Total similarity:
 $\text{sim}(e_i, e_j) = a * \text{sim}_{\text{nbr}}(e_i, e_j) + (1-a) * \text{sim}_{\text{nbr}}(\text{nbr}(e_i), \text{nbr}(e_j))$
 where $\text{nbr}(e)$ denotes the neighborhood (in, out) nodes of e



- ## Maintaining the Order of Comparisons
- In similarity-propagation approaches the order of comparisons is dynamic
 - Graph traversal usually supported by a priority queue (PQ) on the similarity score of nodes
 - As entities are resolved, the PQ is updated for maximizing effectiveness & reducing re-comparisons
 - Different strategies of order maintenance:
 - Based on heuristics
 - type of nodes and edge direction [Dong et al. 2005]
 - degree of nodes [Weis & Naumann 2006]
 - edge weights [Kalashnikov & Mehrotra 2006]
 - Triggered by recent matches [Böhm et al. 2012, Lacoste-Julien et al. 2013]
 - Lazy maintenance [Herschel et al. 2012, Altowim et al. 2014]





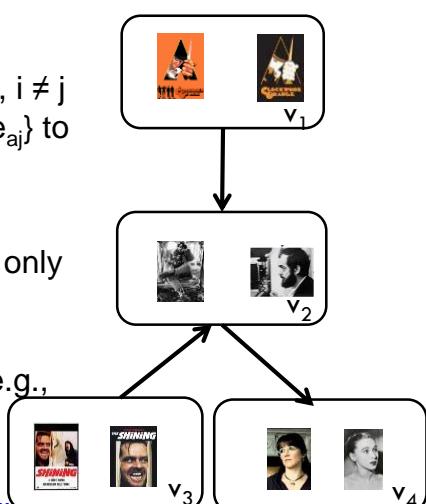
Dependency Graph [Dong et al., 2005]

- Works on an **entity graph** constructed from the relational records
 - **nodes** represent similarity comparisons between pairs of records and/or their attribute values (**real-valued**)
 - **edges** represent match decisions based on the matching of associated nodes (**boolean-valued**)
- A **matching decision** is taken when the real-valued similarity score (between 0 and 1) of a node is above a threshold θ
 - If it exceeds the threshold, it is marked as **match**, otherwise as **undecided**,
 - if no more neighbors are undecided, it is marked as **non-match**
- Idea 1: consider **richer matching evidence**
- Idea 2: **propagate similarity** between matching decisions
- Idea 3: Gradually **enrich references** by merging attribute values (Swoosh-style)



Dependency Graph: Example

- Let E be a set of entity descriptions
 - A **node** $v = \{e_i, e_j\}$, where $e_i, e_j \in E$, $i \neq j$
 - An **edge** $e = (v_a, v_b)$ from $v_a = \{e_{ai}, e_{aj}\}$ to $v_b = \{e_{bi}, e_{bj}\}$ implies $e_{bi}, e_{bj} \in \text{values}(e_{ai}) \cup \text{values}(e_{aj})$
- **Directed edge** when dependency is only in one direction
 - can be stated by domain experts
 - inferred by the data semantics (e.g., keys/foreign keys, rdf properties)
- Include only **nodes whose two entities have the potential to be similar**





Consider Richer Matching Evidence

Fall 2019

[Dong et al., 2005]

- Positive evidence (i.e., constraints for **match** nodes) is captured by the **Boolean similarity** of neighborhood nodes
 - **Strong-boolean**: Resolution implies resolution of neighbour
 - E.g., if two movies are matched then director must also be matched
 - **Weak-boolean**: No direct implication
 - E.g., similarity of two movies increases as their `rdf:labels` are highly similar
- Negative evidence (i.e., constraints for **non-match** nodes) is verified after similarity propagation is performed, and inconsistencies are fixed



137



Fall 2019

Similarity Propagation [Dong et al., 2005]

- Similarity function for node u : $\text{sim}(u) = S_{rv} + S_{sb} + S_{wb}$
 - S_{rv} : from **real-valued** neighbors (decision-tree shape)
 - S_{sb} : from **strong-boolean-valued** neighbors
 - S_{wb} : from **weak-boolean-valued** neighbors
- When a node is matched, the similarity score of its neighbors is re-computed
- Process **converges** if
 - Similarity score is **monotone** in the similarity values of neighbors
 - Compute neighbor similarities only if **similarity increase is not too small**

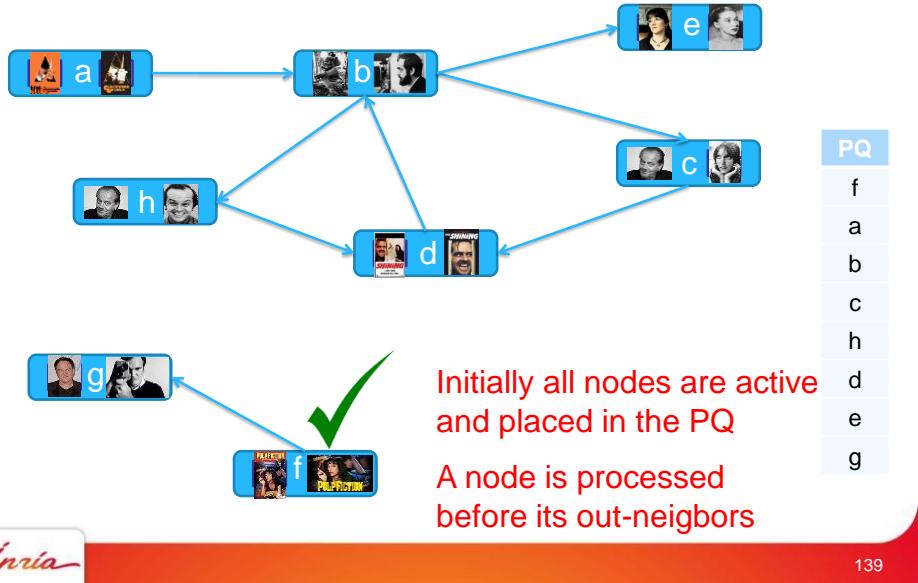


138



Traversing the ER Graph

Fall 2019



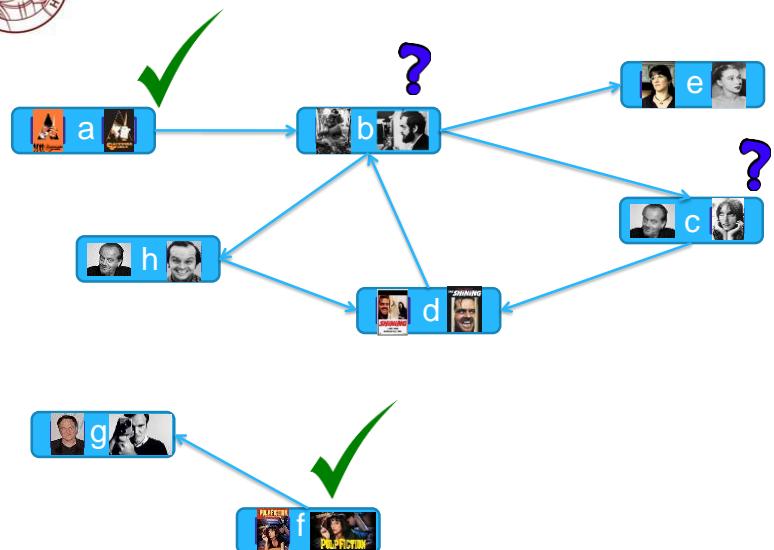
Traversing the ER Graph

Fall 2019





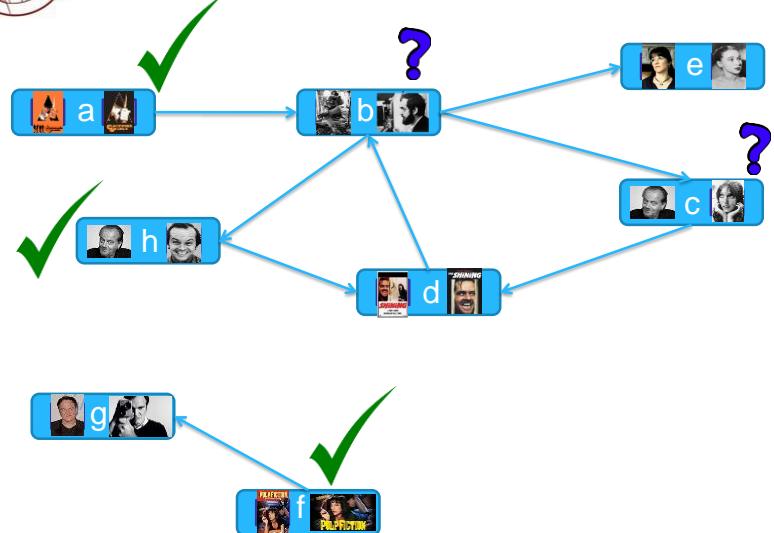
Traversing the ER Graph



141



Traversing the ER Graph

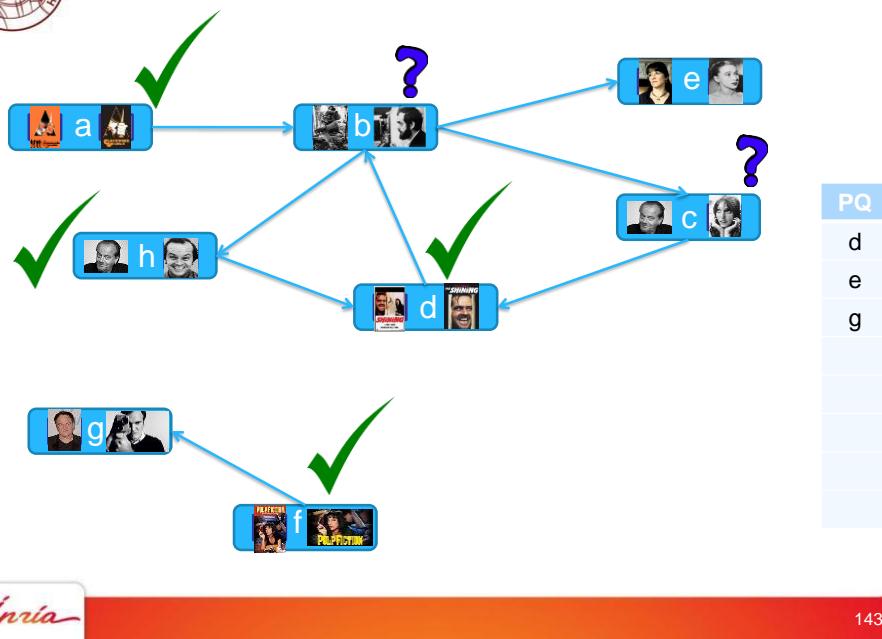


142

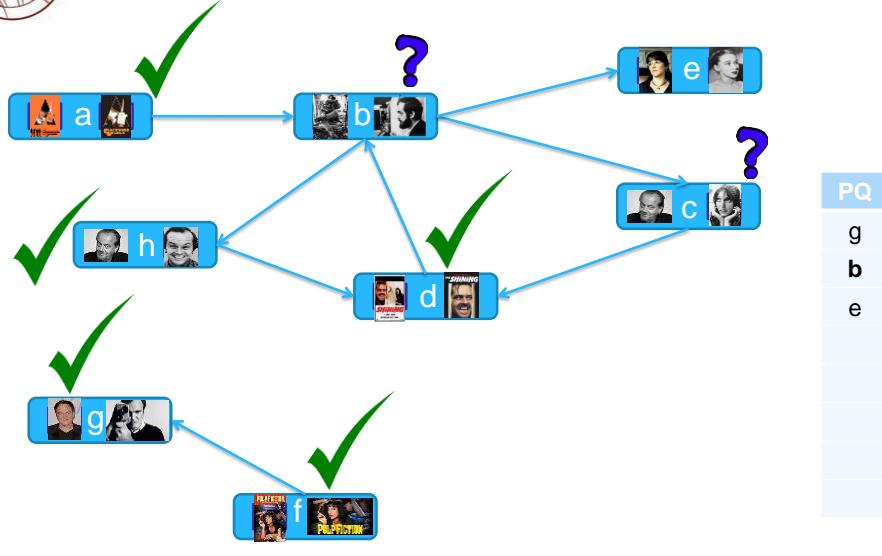




Traversing the ER Graph

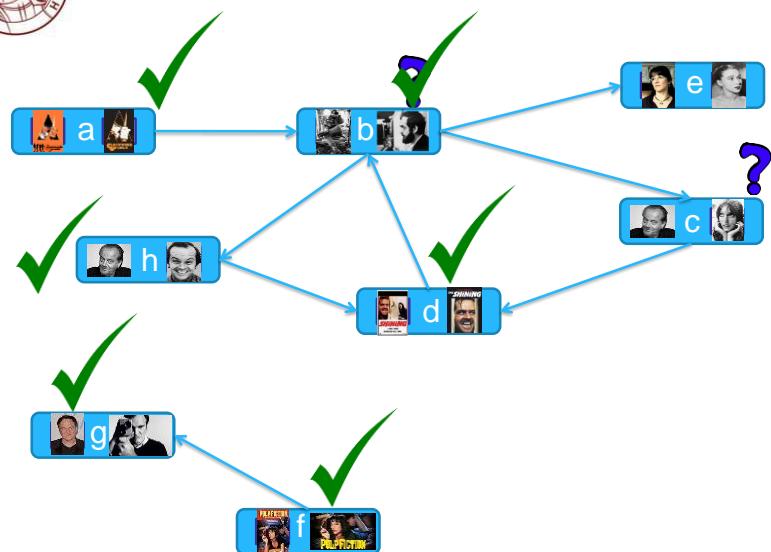


Traversing the ER Graph





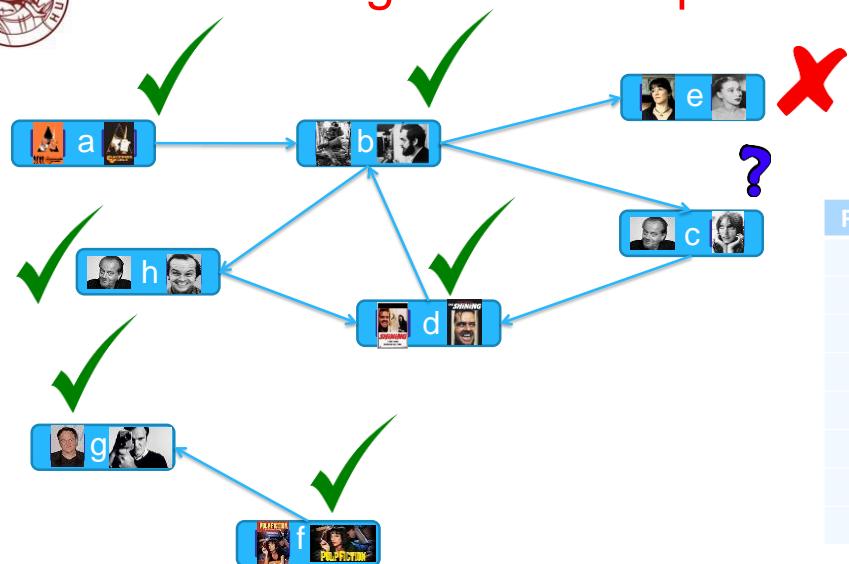
Traversing the ER Graph



145



Traversing the ER Graph

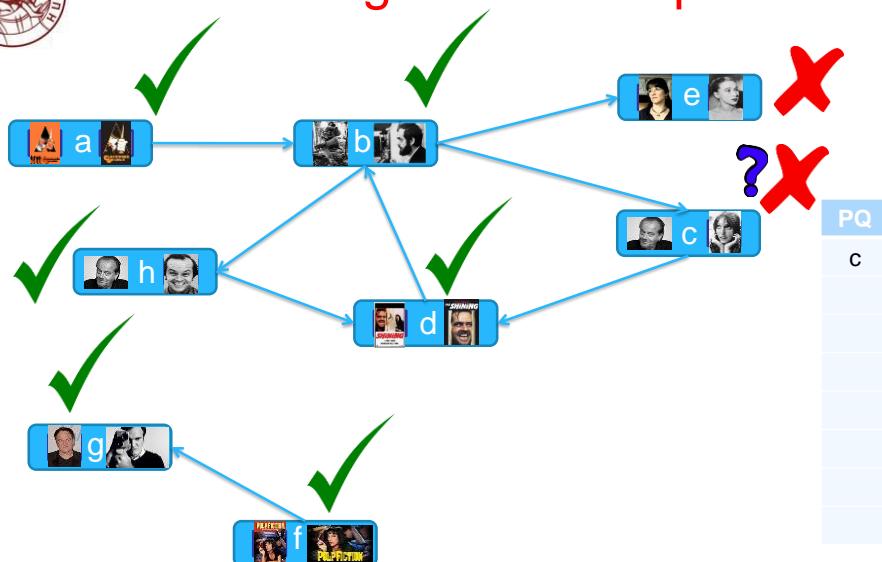


146





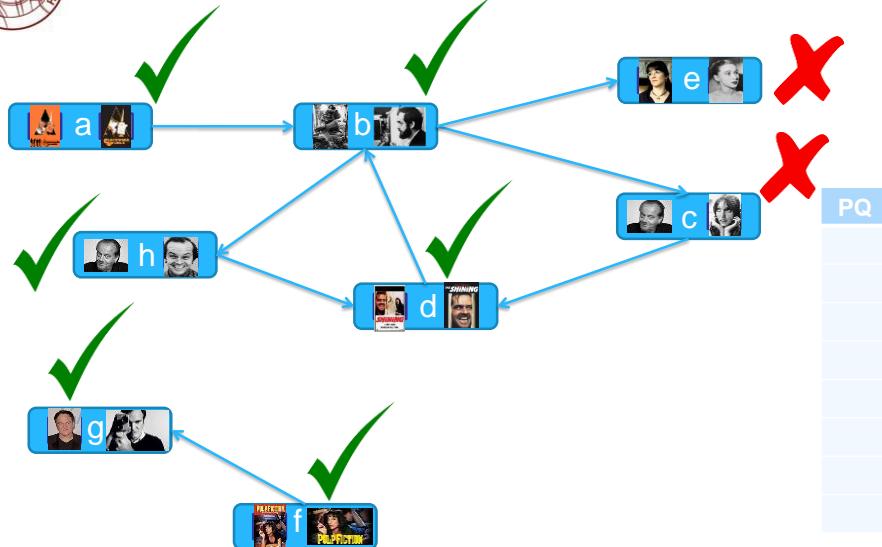
Traversing the ER Graph



147



Traversing the ER Graph



148



Linda [Böhm et al. 2012]

- Works on an entity graph constructed from the RDF descriptions
 - exploits the **unique mapping constraint** between two KBs
- **Key Idea:** the more matching neighbors via similar relationships two descriptions have, the more likely it is that they match
 - **String similarity** of the literal values of entities: checked once
 - **Contextual similarity** of the graph neighbors: checked iteratively
- Two square matrices ($|E| \times |E|$) are used:
 - X captures the **identified matches** (binary values)
 - Y captures the **pair-wise similarities** (real values)
 - Initialization: common neighbors & string similarity of literals
 - Updates: Use the new identified matches of X
- Until the priority queue (extracted from Y) becomes empty:
 - Get the pair (e_i, e_j) with the **highest similarity**: match by default!
 - Update X: matches of e_i are also matches of e_j
 - Update the similarity of nodes influenced by the new matches



149



Linda Algorithm Example

Matches	e1	e2	e3	e4	e5
e1	1	0	0	0	0
e2		1	0	0	0
e3			1	0	0
e4				1	0
e5					1

PQ
e1 – e4
e2 – e4
e1 – e3
e5 – e3
e2 – e3
...

A priority queue, derived by an initial similarity computation between **all pairs**, based on their attribute values



150



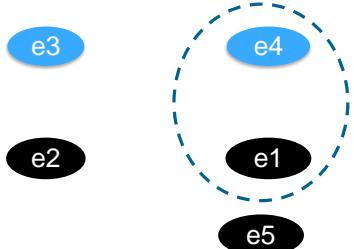
Linda Algorithm Example

Fall 2019

Matches	e1	e2	e3	e4	e5
e1	1	0	0	1	0
e2		1	0	0	0
e3			1	0	0
e4				1	0
e5					1

PQ
e1 – e4
e2 – e4
e1 – e3
e5 – e3
e2 – e3
...

the head of PQ is a match by default



151



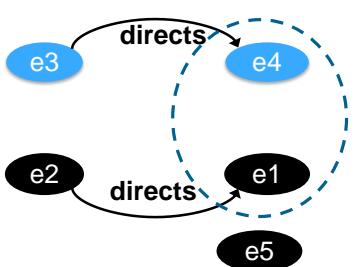
Linda Algorithm Example

Fall 2019

Matches	e1	e2	e3	e4	e5
e1	1	0	0	1	0
e2		1	0	0	0
e3			1	0	0
e4				1	0
e5					1

PQ
e2 – e4
e1 – e3
e2 – e3 ↑
e5 – e3 ↓
...

unique mapping constraint (1-1 Assumption)
similarity re-computation, based on the matching neighbors and the names of the links to them



152

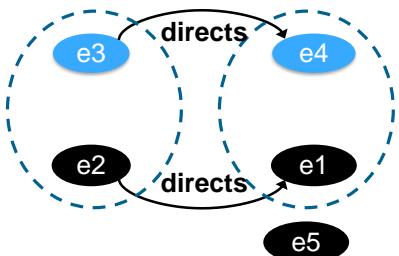


Linda Algorithm Example

Fall 2019

Matches	e1	e2	e3	e4	e5
e1	1	0	0	1	0
e2		1	1	0	0
e3			1	0	0
e4				1	0
e5					1

PQ
e2 – e3
e5 – e3
...



153



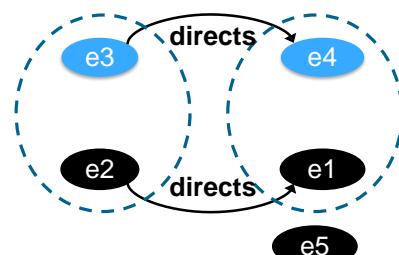
Linda Algorithm Example

Fall 2019

Matches	e1	e2	e3	e4	e5
e1	1	0	0	1	0
e2		1	1	0	0
e3			1	0	0
e4				1	0
e5					1

PQ
e5 – e3
...

unique mapping constraint (1-1 Assumption)



stops when PQ is empty



154

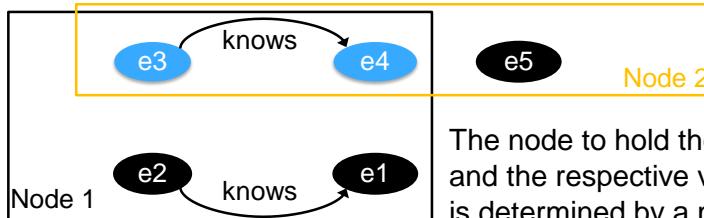


Linda Distributed Version

Fall 2019

Matches	e1	e2	e3	e4	e5	PQ	PQ
e1	1	0	0	0	0	e1 – e4	e5 – e3
e2		1	0	0	0	e2 – e4	e5 – e4
e3			1	0	0	e1 – e3	...
e4				1	0	e2 – e3	Node 2
e5					1	...	

Node 1



The node to hold the PQ entry (a,b) and the respective vertex neighbors is determined by a modulo operation of the first component (a)



155



Fall 2019

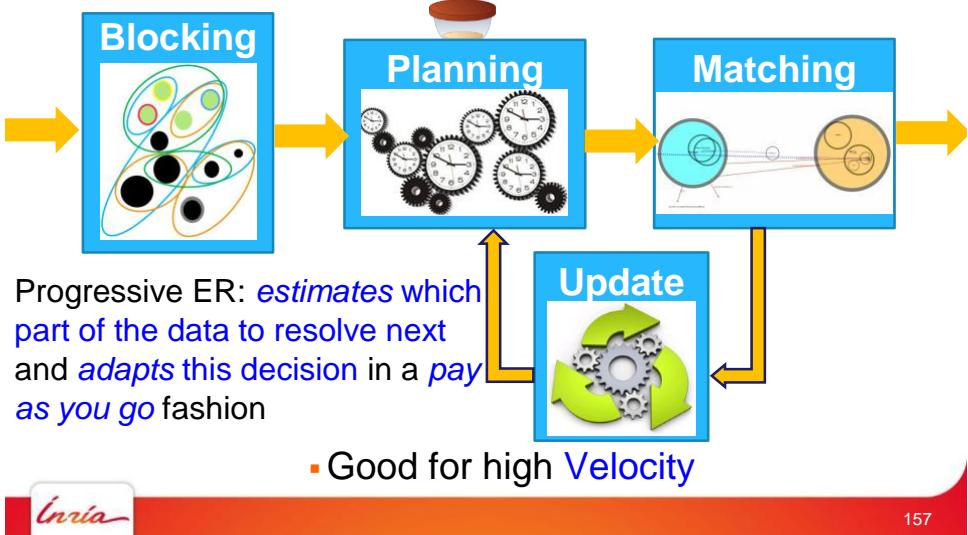
Matching and Resolving Entities (II): Progressive Resolution Techniques





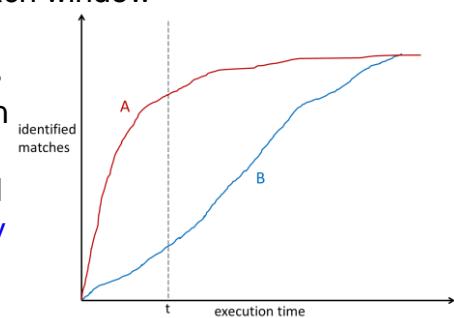
Progressive ER

Optimization: maximize **benefit** (number or type of matches) for a given **cost** (number of comparisons, disk/cloud access)



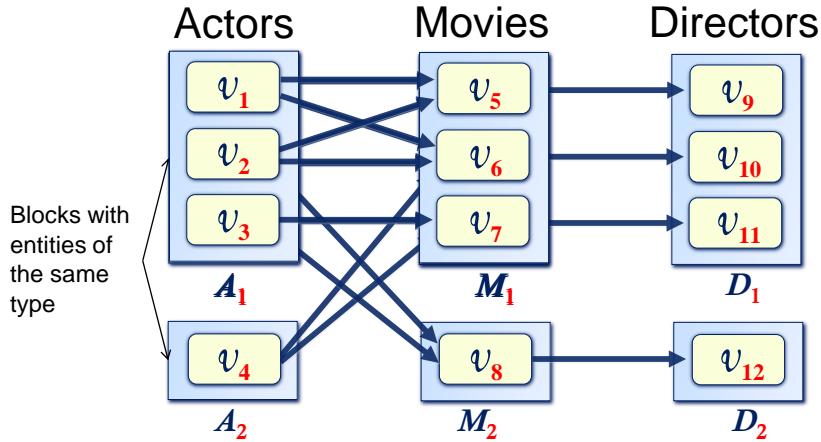
Progressive Approach to Relational Entity Resolution [Altowim et al. 2014]

- Key Idea: divides the ER process into several windows and generates a resolution plan for each window
 - specifies which blocks and entity pairs within these blocks will be resolved during the plan execution phase of a window
 - associates with each identified pair the order in which to apply the similarity functions on the attributes of the two entities
- Lazy resolution strategy to resolve pairs with the smallest cost
 - Unlike single entity type resolution a block based prioritization is significantly more important when resolving multiple types





Relational Progressive ER: Example



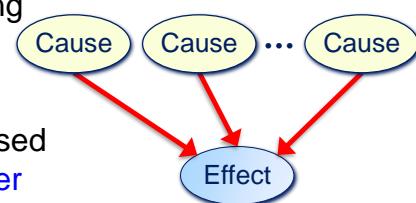
Altowim et al. Progressive Approach to Relational ER VLDB 2014



Matching Probability Estimation

Noisy-OR Model

- The nodes in the ER graph influencing a pair of entities are interpreted as common **causes** to the same **effect**
- The probability that a node is processed at each round depends on the **number of cause nodes** that are matches and/or on the **percentage of matching nodes** in the block of the effect node



Effect: Node considered for matching v_i

Causes: (a) Influencing matching nodes x_j
 (b) Fraction of matching pairs in the block of v_i

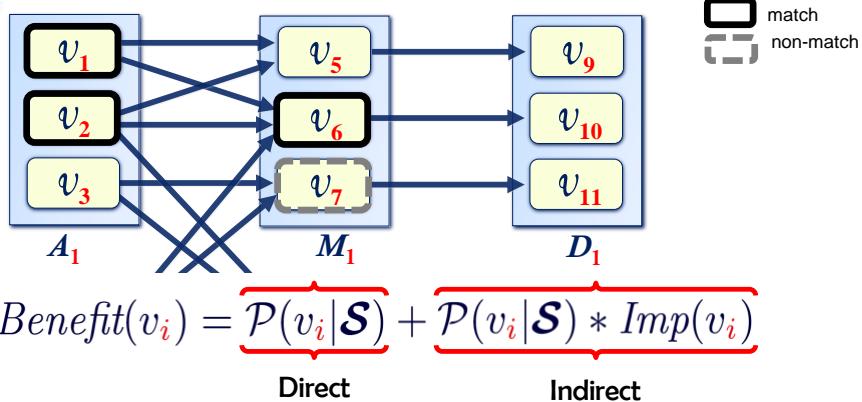
$$\mathcal{P}(v_i) = 1 - \prod_{x_j \in X} 1 - \mathcal{P}(v_i | x_j)$$



Altowim et al. Progressive Approach to Relational ER VLDB 2014



Heuristics for Estimating Node Benefit



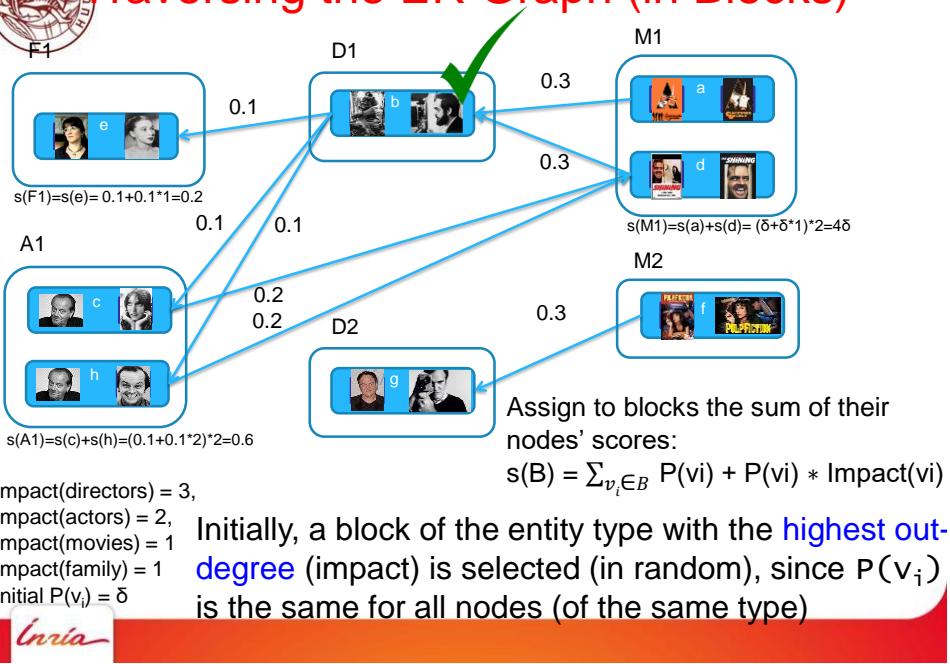
The Impact is estimated for k nodes of the same entity type, as the average number of their direct dependent unresolved nodes

Inria

Altowim et al. Progressive Approach to Relational ER VLDB 2014



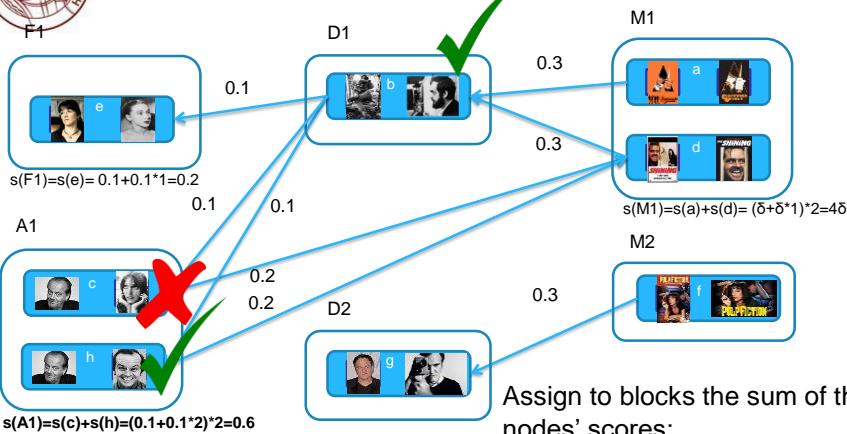
Traversing the ER Graph (in Blocks)



Inria



Traversing the ER Graph (in Blocks)



Impact(directors) = 3,
Impact(actors) = 2,
Impact(movies) = 1
Impact(family) = 1
Initial $P(v_i) = \delta$

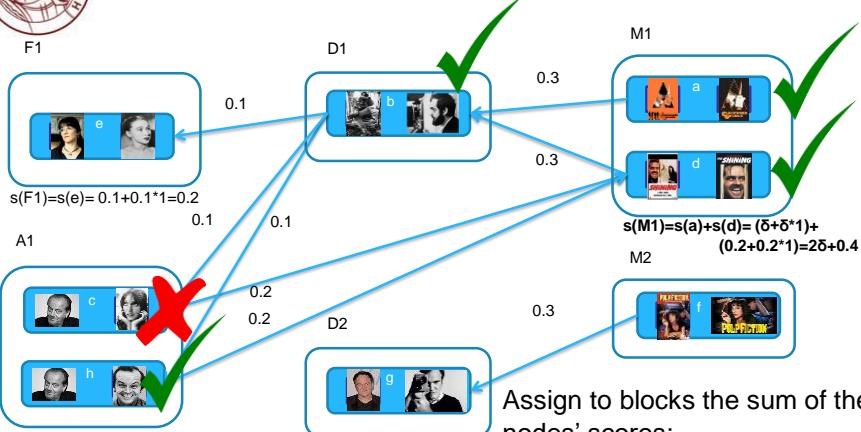
Then, update the scores, **load the block with the highest score** and resolve its nodes

Assign to blocks the sum of their nodes' scores:

$$s(B) = \sum_{v_i \in B} P(v_i) + P(v_i) * \text{Impact}(v_i)$$



Traversing the ER Graph (in Blocks)



Impact(directors) = 3,
Impact(actors) = 2,
Impact(movies) = 1
Impact(family) = 1
Initial $P(v_i) = \delta$

Then, update the scores, **load the block with the highest score** and resolve its nodes

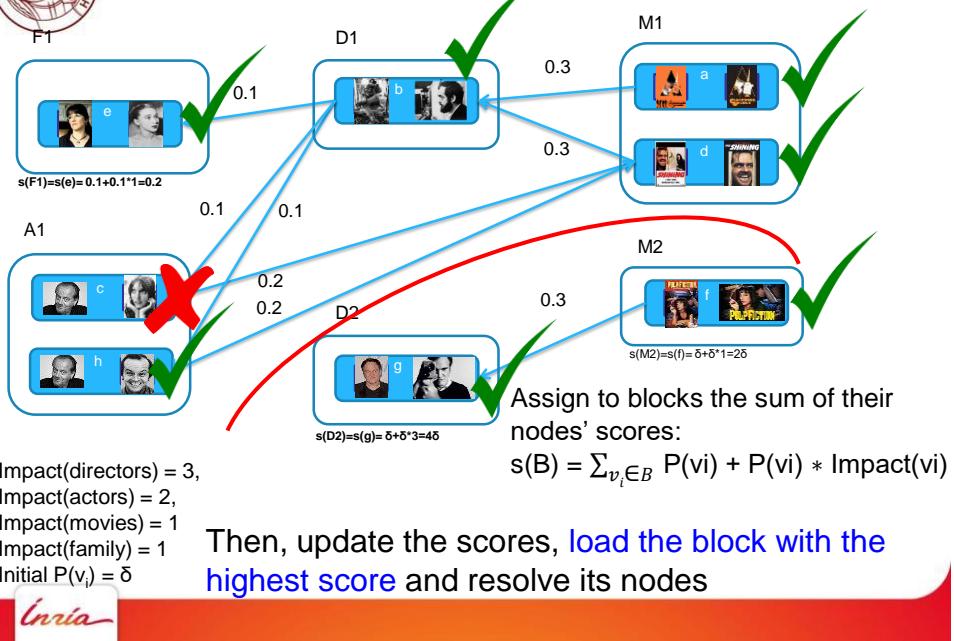
Assign to blocks the sum of their nodes' scores:

$$s(B) = \sum_{v_i \in B} P(v_i) + P(v_i) * \text{Impact}(v_i)$$





Traversing the ER Graph (in Blocks)



Comparison of ER Graph Traversals

Property	[Dong et al. 2005]	[Böhm et al. 2012]	[Altowim et al. 2014]
Seed node selection	any node with no in-edges	the node with the highest content_sim	any node of the entity type with the highest out-degree
Matching condition	similarity is above a threshold value	the highest score (head of PQ)	a (black-box) resolve function returns true
Similarity/Evidence Update (PQ)	weighted sum of in-neighbors similarities	weighted sum of matching in-/out-neighbors similarities	leaky noisy-or of matching in-neighbors
Backtracking	✓	✗	✗
Traversal policy	topological sort	advance in the PQ the neighbors (with most similar relationships) of the recent match	sort PQ based on #in-matches and #out-neighbors





Web-scale Progressive ER: Ongoing Work

Fall 2019

- Relational progressive ER algorithms assume that the **more entity pairs are correctly identified**, the **higher the quality of the result is expected to be**
- We are interested in **characterizing the quality of resolved pairs**
 - # of real-world entities resolved (***shallow strategy***)
 - entity-centric search (*entity coverage*)
 - # of real-world entity graphs resolved (***deep strategy***)
 - entity-centric recommendation (*relationship completeness*)
 - # descriptions resolved for the same real-world entity (***tall strategy***)
 - web-scale knowledge curation (*attribute completeness*)



Fall 2019

Matching and Resolving Entities (II): Challenges and Open Issues





Challenges and Open Issues

- **Tight coupling of Blocking with Iterative Matching/Merging**
 - Better control of block characteristics w.r.t. the entity similarity subsequently used (see [J. Fisher et al. 2015])
- **Progressive ER with Quality Guarantees**
 - guarantees (e.g., coverage) regarding the quality of matches/merges w.r.t. subsequent entity-centric services and data analysis tasks
- **ER for Big Data**
 - Algorithms for high Velocity (see [D. Firmani et al. 2016]), Variety, and Volume entity descriptions (see [Q. Wang et al. 2015] [L. Kolb et al. 2012])
- **Large-Scale ER Testbeds**
 - Real-world ground truth datasets for different match types and open source ER platforms (see [Efthymiou et al. 2015, 2016])

171

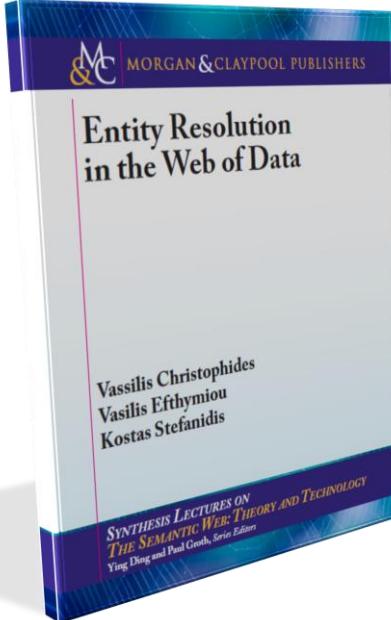


Challenges and Open Issues

- **Crowdsourced ER:**
 - reduce the crowdsourcing cost for obtaining ground truth (see [Chai et al. 2016] [Gokhale et al. 2014] [Wang et al. 2012])
- **Temporal ER:**
 - resolve evolving entity descriptions and analyse the history of descriptions (see [Dong & Tan 2015])
- **Uncertain ER:**
 - consider confidence scores when resolving certain & uncertain entity descriptions (see [Gal 2014] [Demartini et al. 2013])
- **Privacy-aware ER:**
 - Trade-off between entity obfuscation techniques and ER results quality (see [Whang & Garcia-Molina 2013])

172

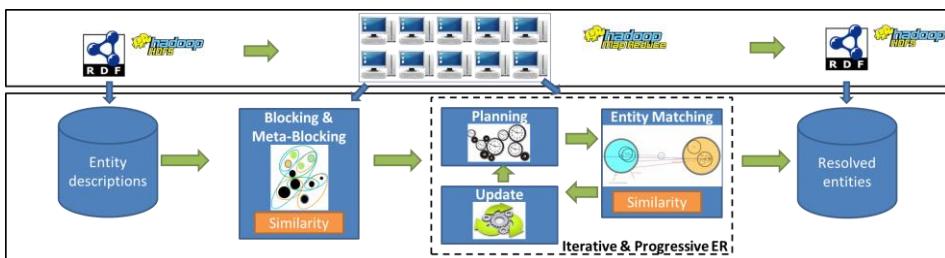




http://www.morganclaypoolpublishers.com/catalog_orig/product_info.php?products_id=823



The Minoan ER Framework



<http://csd.uoc.gr/~vefthym/minoanER> 174



Acknowledgements

- EU FP7-PEOPLE-2013-IRSES **SemData**

- Big Geospatial Data Quality and User Privacy
Défi Mastodons CNRS 2016



License

- These slides are made available under a Creative Commons Attribution-ShareAlike license (CC BY-SA 3.0):
<http://creativecommons.org/licenses/by-sa/3.0/>



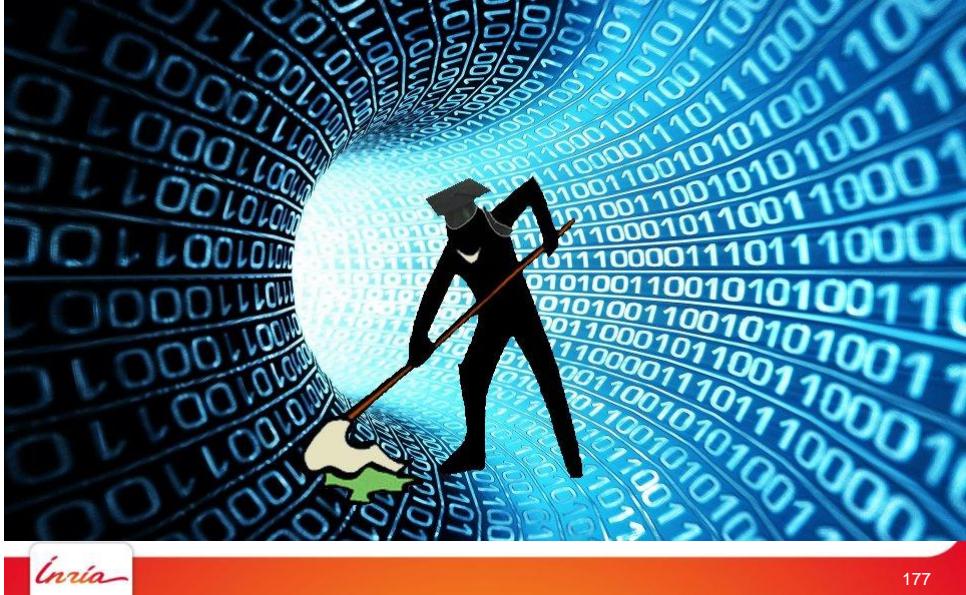
- You can share and remix this work, provided that you keep the attribution to the original authors intact, and that, if you alter, transform, or build upon this work, you may distribute the resulting work only under the same or similar license to this one





Questions?

Fall 2019



177