

Entity Resolution for BIG Data

Blocking-based Entity Resolution in Highly Heterogeneous Information Spaces

George Papadakis

L3S Research Center

Themis Palpanas

University of Trento







Acknowledgements

Ekaterini Ioannou Claudia Niederee Wolfgang Nejdl Peter Fankhauser

dbTrento



Data Evolution

structured precise static (almost)



Figure 1: Entity Relationship Diagram (ERD) for the aer_sight database



dbTrento

Data Evolution

non-structured

uncertain

streaming





large scale data

- very large collections (i.e., terabytes to exabytes)
- scientific, business, user generated

dbTrento



large scale data

streaming data

- sensors, feeds, continuous analytics
- response times in seconds to nanoseconds

dbTrento



large scale data

streaming data

heterogeneous data

- structured, non-structured, text, multimedia
- variety of sources, schemas, representations, models
- computer-generated, human-generated





large scale data

- streaming data
- heterogeneous data
- private data
 - correlating data from multiple sources poses risks
 - credit card history, mobile phone usage, GPS tracking
 - privacy and accountability
 - access granted only to specific person, at specific time, for specific purpose, only for necessary data



large scale data streaming data heterogeneous data private data uncertain data

- imprecision, inconsistencies, incompleteness, ambiguities, latency, deception, approximations, privacy preserving transformations
- process/data/model uncertainty

9



large scale data streaming data heterogeneous data private data uncertain data

confluence of all the above!



Work in dbTrento

large scale data streaming data heterogeneous data private data uncertain data



This Talk

large scale data streaming data heterogeneous data private data uncertain data



This Talk

large scale data (web data sources) streaming data heterogeneous data (no hope for standard schema) private data uncertain data (imprecisions, errors, typos)





Outline

Introduction

- Background
- Standard Blocking for Databases
- Advanced Blocking Methods
 - Q-grams Blocking
 - Suffix Array Blocking
 - The Sorted Neighborhood Approach
 - Canopy Clustering
- Blocking over Dataspaces

Approach

Evaluation



Entities: an invaluable asset

"Entities" is what a large part of our knowledge is about:



dbTrento

However ...

How many names, descriptions or IDs (URIs) are used for the same "entity"?



dbTrento

capital of UK, host city of the IV Olympic Games, host city of the XIV Olympic Games, future host of the XXX Olympic Games, city of the Westminster Abbey, city of the London Eye, the city described by Charles Dickens in his novels, ...

http://sws.geonames.org/2643743/ http://en.wikipedia.org/wiki/London http://dbpedia.org/resource/Category:London



or ...

. . .

How many "entities" have the same name?

0	London, KY		
0	London, Laurel, KY	0	Lon
0	London, OH		261
0	London, Madison, OH		(334)
0	London, AR		(00)
0	London, Pope, AR	0	Lon
0	London, TX		251 Mor
0	London, Kimble, TX		(334)
0	London, MO		X
0	London, MO	0	Lon
0	London, London, MI		122 Van
0	London, London, Monroe, MI		(479
0	London, Uninc Conecuh County, AL		_
0	London, Uninc Conecuh County, Conecuh, AL	0	Lon
0	London, Uninc Shelby County, IN		La
0	London, Uninc Shelby County, Shelby, IN		(858
0	London, Deerfield, WI		
0	London, Deerfield, Dane, WI	0	•••
0	London, Uninc Freeborn County, MN		
0			

dbTrento

idon, Jack 2 Almes Dr ntgomery, AL 4) 272-7005

don, Jack R 1 Winchester Rd ntgomery, AL 36106-3327 4) 272-7005

don, Jack 2 Whitetail Trl 9 Buren, AR 72956-7368 9) 474-4136

don, Jack 0 Vista Del Mar Ave Jolla, CA 92037-4954 8) 456-1850



dbTrento

or ...

How many content types / applications provide valuable information about each of these "entities"?



Pictures and tags about London flickr



Preliminaries

Entity Resolution:

identifies and aggregates the *different* entity profiles/records that actually describe the same real-world object.

Application areas:

Social Networks, census data, price comparison portals, Linked Data

Useful because:

- improves data quality and integrity
- fosters re-use of existing data sources



The entity collections given as input to ER can be of two types:

- **clean**, which are duplicate-free e.g., DBLP, ACM Digital Library, Wikipedia, Freebase,
- **dirty**, which contain duplicate entity profiles in themselves e.g., Google Scholar, Citeseer^X





Types of Entity Resolution

An ER task that receives as input two entity collections can be of the following types:

- **Clean-Clean ER**
- **Dirty-Clean ER** 2.
- **Dirty-Dirty ER** 3.

dbTrento



dbTrento

An ER task that receives as input two entity collections can be of the following types:

- **Clean-Clean ER** (a.k.a. **Record Linkage** in databases) Ι. Given two clean, but overlapping entity collections, identify the entity profiles they have in common. e.g., merge DBLP with ACM Digital Library
- **Dirty-Clean ER** 2.
- **Dirty-Dirty ER** 3.



Types of Entity Resolution

An ER task that receives as input two entity collections can be of the following types:

- **Clean-Clean ER**
- **Dirty-Clean ER** 2.
- **Dirty-Dirty ER** 3.

Identify unique entity profiles contained in union of the input entity collections.

For simplicity, we treat them as equivalent to **Dirty ER** (a.k.a. **Deduplication** in databases), which receives a single dirty entity collection as input.



ER is an inherently quadratic problem (i.e., $O(N^2)$): every entity has to be compared with all others

Does not scale to large entity collections (e.g., Big Data)

Solution: **Blocking**

- similar entities are grouped into blocks
- comparisons are only executed inside blocks





dbTrento

Assumptions:

- Each entity profile corresponds to a single real-world object. Ι.
- Every entity profile consists of a uniquely identified set of name-2. value pairs.
- A-priori known schema 3.
- For each attribute we know some metadata: 4.
 - level of noise (e.g., spelling mistakes, false or missing values)
 - distinctiveness of values
- Two matching profiles are detected if they have at least one 5. block in common.



Standard Blocking for Databases

Simplest form of blocking.

Works as a hash function.

Algorithm:

- Select the most appropriate attribute names w.r.t. noise and distinctiveness.
- Extract a summary from their values in order to form a set of 2. Blocking Key Values (BKVs)
- For each BKV, create one block that contains all entities having 3. this BKV in their transformation.



Example of Standard Blocking



orschungszenth





address=L.A., California, USA

name=William Nicholas

last name=Green

address=L.A., California, USA



Advanced Blocking Methods

Most blocking methods for databases rely on the same signaturebased functionality (review in [Christen, TKDE 2011]).

Main difference (to standard blocking): redundancy \rightarrow overlapping blocks

Drawbacks:

- Schema-dependent Ι.
- Too many parameters to be configured (e.g., which attributes to 2. select, how to combine them in BKVs)

Partial solution:

automatic configuration through machine learning techniques [Bilenko et. al., ICDM 2006] [Michelson et. al., AAAI 2006]



Converts every BKV into a list of q-grams

(i.e., substrings of length q).

For q=2, the keys 91456 and 94520 yield the following blocks:



Drawback:

larger blocks \rightarrow higher computational cost



Suffix Array Blocking [Aizawa et. al., WIRI 2005][de Vries et. al., CIKM 2009]

Converts every BKV to its list of suffixes that are longer than a predetermined minimum length I_{min}.

For I_{min} =3, the keys 91456 and 94520 yield the blocks:



Drawback:

larger blocks \rightarrow higher computational cost



Sorted Neighborhood [Hernandez et. al., SIGMOD 1995]

 Entities are sorted in ⁹¹⁴⁵⁶ alphabetic order of BKVs.
A window of fixed size ⁹⁴⁵²⁰ slides over the sorted list.





Sorted Neighborhood [Hernandez et. al., SIGMOD 1995]

Ι.	Entities are sorted in	91456
	alphabetic order of BKVs.	
2.	A window of fixed size	94520
	slides over the sorted list.	





Sorted Neighborhood [Hernandez et. al., SIGMOD 1995]

91456

94520

- Entities are sorted in alphabetic order of BKVs.
- A window of fixed size 2. slides over the sorted list.

Drawback: the fixed window may be

- too narrow in some cases, leaving out matches,
- too wide in other cases, including many irrelevant entities.

Solution:

adaptable window size based on the similarity of BKVs (Adaptive Sorted Neighborhood [Yan et.Al., JCDL 2007])





Canopy Clustering [McCallum et. al., KDD 2000]

Not based on signatures.

Instead, entity profiles are compared with a cheap string similarity metric *m* (e.g., TF-IDF, Jaccard similarity).

Algorithm:

- Select a random entity profile e from the pool of entities
- Create a new block containing 2. all e_i s.t. $m(e_1,e_i) < t_1$
- Remove from the pool of entities 3. all e_i s.t. $m(e_1, e_i) < t_2 (t_2 < t_1)$
 - Repeat until the pool is empty



4.



They include Web 2.0 data, Semantic Web, Dataspaces.

Voluminous, (semi-)structured datasets.

- DBPedia 3.4: 36.5 million triples and 2.1 million entities
- ▶ BTC09: 1.15 billion triples, 182 million entities.

Users are free to insert not only attribute values but also attribute names \rightarrow high levels of heterogeneity.

- DBPedia 3.4: 50,000 attribute names
- Google Base: 100,000 schemata and 10,000 entity types BTC09: 136K attribute names

Large portion of data originating from automatic information extraction techniques \rightarrow noise, tag-style values





Existing blocking techniques are inapplicable, due to:

- Loose schema binding.
- High levels of heterogeneity.
- Noise and missing attribute names or values.

These settings call for blocking methods that:

- Are robust to noise.
- Are decoupled from schema information.
Exemplary Challenges of Big Data

orschungszenth





Problem Definition

Metrics for assessing block quality:

- Pair Completeness:
- **Reduction Ratio:**

detected matches PC = existing matches

 $RR = 1 - \frac{\text{method} comparisons in \# of}{RR}$ baseline comparisons comparisons)

Problem Definition:

Given two duplicate-free entity collections (Clean-Clean ER), cluster their entities into blocks and process them so that both PC and RR are maximized.

The same applies to Dirty ER.

disclaimer:

Precision of entity matching is dependent on the entity similarity measures, and is orthogonal to the above problem.



Outline

- Introduction
- Approach
 - Core Ideas
 - Effectiveness Layer
 - Type of pair-wise comparisons
 - Metric Space for Blocking Methods
 - Meta-blocking Layer
 - Efficiency Layer
- Evaluation



Framework Outline

Three-layered approach*:

Effectiveness Layer



Build blocks in a way that involves:

- Attribute-agnostic functionality
- Redundancy for robustness and high effectiveness

Meta-blocking 2.

Restructure a block collection into a new one with fewer comparisons, but equally high effectiveness.

Efficiency Layer 2.

Process blocks s.t. entire blocks or individual comparisons that do not contribute to effectiveness are discarded (thus increasing efficiency).



*Code and data are fully available at: sourceforge.net/projects/erframework/.



dbTrentc

Effectiveness Layer

Goal:

place every pair of matching entities in at least one common block. Solution:

- redundancy to reduce the likelihood of missed matches,
- attribute-agnostic functionality to be able to handle Big Data. Drawbacks:
- the blocks are overlapping (i.e., repeated comparisons),
- high number of comparisons between irrelevant entities.

Three families of approaches:

- Token Blocking
- Attribute Clustering (Clean-Clean ER) 2.
- URI Semantics Blocking (Dirty ER) 3.



Attribute-agnostic blocking scheme:

- completely ignores attribute names.
- considers all attribute values.

Functionality:

- given an entity profile, it transforms it into the set of tokens contained in its values.
- creates one block for each distinct token \rightarrow each block 2. contains all entities with the corresponding token*.

Redundancy is inherent!



*Each block should contain at least two entities.



Token Blocking Example



Entity 2

Entity 4

d

Entity 2

Entity 4



Attribute-Clustering Blocking [Papadakis et. al., TKDE2013]

Goal

group attribute names into clusters s.t. we can apply Token Blocking independently in each cluster, yielding more efficient but equally effective blocks

Algorithm

Input data: N_1 , N_2

Parameters: a string similarity metric m, a representation model r**Procedure:**

for every attribute name n_i in N_1

find the most similar n_i in N_2 s.t. $m(r(n_i), r(n_j)) > 0$

connect them with an edge

do the same for every attribute name n_i in N_2

get transitive closure

get connected components

merge together singleton clusters



Versatile settings:

- character n-grams with Jaccard similarity Ι.
- token n-grams with TF-IDF 2.
- n-gram graphs with value similarity (graph metric) 3.

The last one performs better, but computationally expensive.

Similar to Schema Matching, but fundamentally different:

- Associated attribute names do not have to be semantically equivalent. They only have to produce good blocks.
- All singleton attributes are associated with each other. 2.
- Schema matching approaches do not scale to the highly heterogeneous big 3. data.

45



Evidence for URI Semantics Blocking

Attribute-agnostic blocking leads to high levels of redundancy. For Semantic Web data, three sources of evidence create blocks of lower redundancy:

1. Semantics in entity URIs (i.e., Infix) [Papadakis et al., iiWAS 2010] Infixes may consist of several concatenated tokens.

> Prefix http://db1p.13s.de/d2r/resource/publications/books/sp/wooldridgeV99 http://bibsonomy.org/uri/bibtexkey/books/sp/wooldridgeV99

- 2. Relations between entities (i.e., Infix Profile) The URIs of attribute values can be represented by their Infixes.
- 3. Literal values (i.e., *Literal Profile*)

Values th

at do not cor	respond to blank nodes or URIs can	be represented by their tokens.
URL: birthname: dateOfBirth:	<http: barack_obama="" dbpedia.org="" resource=""> "Barack Hussein Obama II" "1961-08-04"</http:>	Barack_Obama Michelle_Obama Joe_Biden ^{Hawaii}
birthPlace: shortDescription: spouse: Vicepresident:	"Hawaii" <http: dbpedia.org="" hawaii="" resource=""> "44th President of the United States of America" <http: dbpedia.org="" michelle_obama="" resource=""> <http: dbpedia.org="" joe_biden="" resource=""></http:></http:></http:>	Barack 08 America States 01 Obama 04 20 44th 2009 of Hussein Hawaii United 1961 the II President





URI Semantics Blocking [Papadakis et al., WSDM2012]

The above sources of evidence lead to 3 blocking methods:

I. Infix Blocking

every block contains all entities whose URI has a specific Infix

2. Infix Profile Blocking

every block corresponds to a specific Infix (of an attribute value) and contains all entities having it in their Infix Profile

3. Literal Profile Blocking

every block corresponds to a specific token and contains all entities having it in their Literal Profile

Individually, these methods may not cover all entities (e.g., Infix Blocking does not cover blank nodes). However, they are complementary and can be combined into composite blocking methods for higher effectiveness.



Outline

- Introduction
- Approach
 - Core Ideas
 - Effectiveness Layer
 - Type of pair-wise comparisons
 - Metric Space for Blocking Methods
 - Meta-blocking Layer
 - Efficiency Layer
- Evaluation



Meta-blocking Layer [Papadakis et. al., TKDE]

Main idea

block assignments provide valuable evidence for the similarity of entities: the more blocks two entities share, the more similar and more likely they are to be matching

Goal:

restructure a given block collection into a new one that contains substantially lower number of comparisons (RR»0), while being equally effective (same PC).

dbTrento



- Every comparison between entity profiles p_i and p_i belongs to one of the following types:
- Matching comparison if $P_i \equiv P_i$. 1.
- Redundant comparison if p_i and p_i have been compared in a 2. previously examined block.
- Superfluous comparison if p_i or p_i or both of them have been 3. matched to some other entity (Clean-Clean ER).
- Non-matching comparison if $P_i \neq P_i$. 4.
- To enhance efficiency without any impact on effectiveness, we should discard the last three types of comparisons.



Main idea

block assignments provide valuable evidence for the similarity of entities: the more blocks two entities share, the more similar and more likely they are to be matching

Goal:

restructure a given block collection into a new one that contains substantially lower number of redundant and non-matching comparisons ($RR \gg 0$), while being equally effective (same PC).

Solution: **Blocking graph**

- Nodes \rightarrow entities
- Edges \rightarrow between entities co-occurring in at least one block

Weights \rightarrow how similar are the adjacent entities





The graph pruning algorithms requires setting a threshold. Ideally, we could set pruning thresholds using PC and RR.

However, PC and RR can only be measured *a-posteriori*: have to first construct and examine all blocks

Instead, we need an *a-priori* estimation of the PC and RR values: in order to guide the block restructuring process without executing any comparisons



BC-CC metric space:



Blocking Cardinality (BC)

average block assignments per entity: highly correlated with PC

Comparisons Cardinality (CC)

average block assignments per comparison: highly correlated with RR



Block Enhancement on BC-CC Space

Enhancing Effectiveness

combining complementary atomic methods into composite ones



dbTrento





Enhancing Effectiveness (e.g., Method₁, Method₂ \rightarrow Method₃)

combining complementary atomic methods into composite ones







Enhancing Effectiveness (e.g., Method₁, Method₂ \rightarrow Method₃)

combining complementary atomic methods into composite ones

Enhancing Efficiency

- Meta-blocking
- **Block Processing**







Enhancing Effectiveness (e.g., Method₁, Method₂ \rightarrow Method₃)

combining complementary atomic methods into composite ones

Enhancing Efficiency (e.g., Method₃ \rightarrow Method₄)

- Meta-blocking
- **Block Processing**





Transforming the pruned blocking graph into a block collection.

For undirected blocking graphs:

every retained edge creates a block of minimum size

For directed blocking graphs:

for every node (with retained outgoing edges), we create a new block containing the corresponding entities





Outline

- Introduction
- Approach
 - Core Ideas
 - Effectiveness Layer
 - Type of pair-wise comparisons
 - Metric Space for Blocking Methods
 - Meta-Blocking Layer
 - Efficiency Layer
- Evaluation





Efficiency Layer

Goals:

- I. eliminate repeated comparisons,
- 2. discard superfluous comparisons,
- 3. avoid non-matching comparisons.

without affecting effectiveness.

Techniques:

Comparison's Type G r Superfluity Repeat Non-ma а Method Method metho n u 1. Block Pu **Block**refinement 2. Block Pri a r Comparis **Comparison-**Duplicate Comparison t refinement Propagation Propagation Prunin

dbTrento

tch	Scheduling
od	method
irging	Block
uning	Scheduling
son	Comparison
g	Scheduling



Block Purging [Papadakis et al., WSDM2011] & [Papadakis et al., WSDM2012]

Oversized blocks: large number of comparisons, the vast majority of which is redundant, non-matching and superfluous.

Block Purging: discards oversized blocks by setting an upper limit

on size of each block [Papadakis et al., WSDM 2011], or

on number of comparisons it contains [Papadakis et al., WSDM 2012]











Purged Set of Blocks



Duplicate Propagation [Papadakis et al., WSDM2011]

In Clean-Clean ER, matched Entities do not need to be compared again. They can be propagated through a central data structure that contains the matched duplicates from each dataset.



```
Duplicates Dataset 1: {null}
Duplicates Dataset 2: {null}
```

```
Duplicates Dataset 1: {Entity 1}
Duplicates Dataset 2: {Entity 3}
```

```
Total Comparisons: 3
Duplicates Dataset 1: {Entity 1, Entity 2}
Duplicates Dataset 2: {Entity 3, Entity 4}
```

```
Total Comparisons: 3!
Duplicates Dataset 1: {Entity 1, Entity 2}
Duplicates Dataset 2: {Entity 3, Entity 4}
```



Enhance effect of Duplicate Propagation, by detecting duplicates as early as possible:

order blocks so as to examine first those blocks that are highly likely to contain duplicates

This likelihood is expressed through the utility measure that is assigned to each block:

$$U_{i} = \frac{gain_{i}}{cost_{i}} = \frac{1}{max(|b_{i,1}|, |b_{i,2}|)}$$

where $|b_{i,1}|$ and $|b_{i,2}|$ are number of entities from dataset 1 and 2, that are contained in block b_i



Block Scheduling Example

Initial Set of Blocks



Total Comparisons: 3 (see slide 15)





Total Comparisons: 0
Duplicates Dataset 1: {null}
Duplicates Dataset 2: {null}

$$U_{Veneman}=1$$

Total Comparisons: 1
Duplicates Dataset 1: {E2}
Duplicates Dataset 2: {E4}
 $U_{Ann}=1$
Total Comparisons: 1
Duplicates Dataset 1: {E2}
Duplicates Dataset 2: {E4}
 $U_{unicef}=0.5$

Total Comparisons: 2 Duplicates Dataset 1: {E1, E2} Duplicates Dataset 2: {E3, E4}



Effect of Block Scheduling

cost of identifying **new** duplicates is called duplicate overhead (dh): $dh_k = \frac{comparisons_k}{duplicates_k}$

comparisons_k: # of comparisons executed after examining the (k-I)-th block containing new duplicates

 $duplicates_k$: # of matched pairs of entities in the k-th block containing new duplicates.

dh *increases* as we move from first block to next ones:

due to Block Scheduling, the lower the execution order of a block is, the more comparisons and the less unique duplicates 4th it entails





dh provides a good estimation of the duplicates that remain to be detected \rightarrow the higher its value is, the more unlikely it is to identify new duplicates in the remaining blocks at a reasonable cost (i.e., #comparisons).

Therefore, the ER process can be terminated prematurely, when dh gets high.

Block Pruning sets an upper limit on duplicate overhead, called dh_{max}, based on the total number of comparisons n contained in the original set of blocks: $dh_{max} = 10^{\log n/2}$

As soon as dh exceeds dh_{max} , the entire process is terminated.



Goal : eliminate all redundant comparisons.

Naïve solution: store all executed comparisons in a hash table \rightarrow not scalable

Our approach: indirect propagation of all executed comparisons in three steps

- **Block Enumeration**: assign to each block an index denoting its position on the block processing list.
- **Inverted Index**: Data structure associating entities with block indices. 2.



Two entities are compared iff they satisfy the **Least Common Block** 3. **Index Condition**: the index of the current block is their minimum common block index. E.g., p_1 and p_3 are compared in b_1 , but the condition does not hold in b_{5} .



Goal: discard non-matching comparisons.

How can we infer that two entities are non-matching without actually comparing them?

Examine number of blocks they have in common! Comparison Pruning discards a comparison between p_i and p_i if it does not satisfy the following condition:

$$JS_{min} \leq ES(p_i, p_j)$$

where: $JS(p_i, p_i)$ is the Jaccard coefficient of their block lists, and $a \cdot min(iBC_{t,c}^{\varepsilon_1}, iBC_{t,c}^{\varepsilon_2})$

 $ES_{min} = \frac{a \cdot min(a \in i,c)}{iBC_{t,c}^{\mathcal{E}_1} + iBC_{t,c}^{\mathcal{E}_2} - a \cdot min(iBC_{t,c}^{\mathcal{E}_1}, iBC_{t,c}^{\mathcal{E}_2})}$





Comparison Scheduling [Papadakis et. al., TKDE1]

Enhances the effect of Duplicate Propagation: it schedules the execution of individual comparisons so that those involving real matches are executed first.

The utility of every comparison between p_i and p_i is: $u(p_i,p_i) = JS(p_i,p_i) \cdot ICF(p_i) \cdot ICF(p_i)$

where $JS(p_i, p_i)$ is the Jaccard coefficient of the block lists, and ICF(p_k) is the inverse comparison frequency of p_k .

To reduce its high computational cost, it is usually applied after Comparison Propagation and Comparison Pruning.



Composing ER workflows [Papadakis et. al., TKDE2013]



1 st step	2 nd step	3 rd step	4 th step	
Block Building	Meta- blocking	Core Methods	Scheduling Methods	
Token Blocking OR Agnostic Clustering OR OR Total Description	WEP OR CEP OR WNP OR OR CNP	Block Purging AND Duplicate Propagation	Block Scheduling OR Comparison Scheduling	Co

Actual selection depends on the application requirements and the available resources.

In general, comparison-refinement methods are more accurate,

but more time- and resource-consuming than block-refinement.





Outline

- Introduction
- Approach
- Evaluation
 - Experimental Settings
 - Block Building Performance
 - Scalability of Token Blocking
 - Block Purging Performance
 - Meta-blocking Performance
 - Block Processing Performance


Experimental Settings

Metrics

Pair Completeness

Reduction Ratio:

$$PC = \frac{detected_matches_ecall}{existing_matches}$$
$$RR = 1 - \frac{method_comp(anisions)}{baseline_comparisons}$$

Da	atasets			
		BT		(
		(Dir	ty EK)	30
	Entities	182	million	1,19
	Name-Value Pairs	I.15	billion	17,4
	Duplicates	,59 (IFP)	5,988,554 (Same-As)	



Note I: out of 43,75 million distinct triples of D_{DBPedia}, only 10,36 million (<25%) are common.

Note 2: IFP corresponds to pairs of matching entities that are **implicitly** denoted via their common values for some Inverse Functional Properties. Same-As corresponds to pairs of matching entities that are explicitly denoted through same-as statements.

s in # of comparisons) S





We considered subsets of $D_{DBpedia}$.

For each subset we estimated the following metric:

average comparisons comparisons per entity = DBPedia_{3.0rc} | + | DBPedia_{3.4} |



3-4.5 orders of magnitude less comparisons (efficiency) PC always above 95% (recall)

Our Approach Naive Method [1190734, 2164058]



Block Building Performance

DBPedia	Comparisons	Duplicates	PC	RR
Token Blocking	6.18 · 10 ¹²	892,560	99.99%	_
Term Vector EM	6.38 · 10 ¹²	891,546	99.99%	_
Term Vector AC	6.18 · 10 ¹²	891,560	99.99%	0.01%
Trigrams AC	1.05 · 10 ¹²	892,425	99.98%	83.04%
Trigram Graphs AC	1.03 · 10 ¹²	892,516	99.99%	83.39%

BTC09	Comparisons	PC IFP	PC SameAs	RR
Token Blocking	2.59 · 10 ¹⁶	99.32%	92.26%	_
Infix Blocking	6.24 · 10 ¹²	59.31%	49.60%	99.98%
Infix Profile Blocking	1.07 · 10 ¹⁵	84.19%	99,99%	95.86%
Literal Profile Blocking	2.67 · 10 ¹⁵	94.49%	24.51%	89.68%
Infix+Infix Profile	1.10 · 10 ¹⁵	85.16%	87.14%	95.77%
Infix+Literal Profile	3.25 · 10 ¹⁵	96.43%	65.67%	87.46%
Infix Profile+Literal Profile	3.81 · 10 ¹⁵	96.64%	52.09%	85.30%
All (Combined)	3.90 · 10 ¹⁵	97.98%	91.13%	84.96%
σ				75



Block Purging Performance

DBPedia	Comparisons	Duplicates	PC	RR
Token Blocking	5.68 · 10 ¹⁰	891,767	99.91%	99.08%
Term Vector EM	4.34 · 10 ¹⁰	891,709	99.90%	99.32%
Term Vector AC	5.68 · 10 ¹⁰	891,767	99.91%	99.08%
Trigrams AC	3.06· 10 ¹⁰	892,402	99.98%	97.08%
Trigram Graphs AC	2.42 · 10 ¹⁰	892,463	99.99%	97.64%

	BTC09	Comparisons	PC IFP	PC SameAs	RR
	Token Blocking	15.49 · 10 ¹¹	96.79%	60.52%	99.99%
	Infix Blocking	$0.004 \cdot 10^{11}$	58.85%	49.54%	99.99%
	Infix Profile Blocking	1.62 · 10 ¹¹	76.53%	41.36%	99.99%
	Literal Profile Blocking	9.37 · 10 ¹¹	94.34%	18.29%	99.99%
	Infix+Infix Profile	1.69 · 10 ¹¹	72.64%	86.60%	99.99%
	Infix+Literal Profile	8.79 · 10 ¹¹	94.48%	63.75%	99.99%
	Infix Profile+Literal Profile	12.00 · 10 ¹¹	93.57%	50.03%	99.99%
	All (Combined)	11.64 · 10 ¹¹	95.37%	89.35%	99.99%
F					76





Block Processing Performance

		Comparisons	RR	Duplicates	ΡС	Time
	Block Purging	2.42 · 10 ¹⁰	-	892,463	99.99%	0.05
	Block Scheduling	1.55 · 10 ¹⁰	93.58%	892 <i>,</i> 463	99.99%	0.16
vvr ₁	Block Pruning	7.24 · 10 ⁷	99.70%	879,446	98.53%	0.01
	Comparison Propagation	1.24 · 10 ¹⁰	48.86%	892,463	99.99%	5.75
WF_2	Block Scheduling	9.20 · 10 ⁸	96.20%	892 <i>,</i> 463	99.99%	0.16
	Comparison Pruning	4.98 · 10 ⁷	99.79%	837,286	93.80%	4.14
	Comparison Propagation	1.24 · 10 ¹⁰	48.86%	892,463	99.99%	5.75
WF_3	Comparison Pruning	4.32 · 10 ⁸	98.21%	837,286	93.80%	4.14
	Comparison Scheduling	4.46 · 10 ⁷	99.82%	837,286	93.80%	0.51



dbTrento

Meta-blocking Performance

		Comparisons	RR	ΡС	ΔΡC			
	Token Blocking	3.98 · 1010	-	99.91%	-			
WEP	ARCS	2.85 · 10 ⁸	99.28%	92.45%	-7.45%			
	CBS	3.40 · 10 ⁹	91.46%	95.47%	-4.42%			
	ECBS	5.77 · 10 ⁹	85.50%	99.66%	-0.23%			
	JS	1.11 · 10 ¹⁰	71.80%	99.73%	-0.16%			
	EJS	1.10 · 10 ¹⁰	72.32%	99.77%	-0.11%			
	ARCS	1.85 · 10 ⁹	95.34%	99.41%	-0.48%			
	CBS	3.57 · 10 ⁹	91.04%	99.35%	-0.54%			
CEP	ECBS	9.94 · 10 ⁹	75.02%	99.75%	-0.14%			
	JS	1.96 · 10 ¹⁰	50.76%	99.87%	-0.02%			
	EJS	1.99 · 10 ¹⁰	49.74%	99.88%	-0.01%			

	Comp.	RR	PC					
			ARCS	CBS	ECBS	JS	EJS	
WNP	0.26 · 10 ⁸	99.94%	79.46%	51.71%	61.14%	82.09%	79.61%	
CNP	0.50 · 10 ⁸	99.88%	93.43%	92.35%	94.05%	95.57%	95.99%	



Meta-blocking Time Requirements

Materialization Time (*MT*): time for Graph Building and Edge Weighting.

Restructure Time (*RT*): time for Graph Pruning and Block Collecting.

Comparisons Time (*CT*): time for executing the retained comparisons.

Performance over DBPedia in hours, using Intel Xeon E5472 3.0 GHz and 16GB of RAM. Profile comparison was done with Jaccard similarity.

		_			
		MT	RT	СТ	Total
Token Bl.		0	0	142.37	142.37
	ARCS	7.48	17.00	25.13	49.62
W	CBS	8.08	8.53	26.81	43.42
Ĕ	ECBS	11.69	9.40	27.93	49.03
Р	JS	7.74	6.75	26.76	41.26
	EJS	6.87	7.12	32.46	46.45
	ARCS	6.87	16.08	1.81	24.76
С	CBS	7.12	6.97	4.64	18.73
Ĕ	ECBS	7.26	7.40	2.84	17.49
Р	JS	6.81	6.56	0.28	13.66
	EJS	6.81	7.45	0.25	14.51
	ARCS	6.81	17.58	13.05	40.64
W	CBS	6.84	6.87	29.08	49.99
Ň	ECBS	6.84	7.08	35.56	56.84
Ρ	JS	7.36	9.07	26.74	42.24
	EJS	6.81	7.44	36.85	58.83
	ARCS	6.81	17.58	3.49	27.89
С	CBS	6.84	6.87	5.35	19.06
Ň	ECBS	6.84	7.08	2.05	15.97
Р	JS	7.36	9.07	0.78	17.21
	EJS	6.81	7.44	0.85	15.10



Conclusions

Entity resolution (ER) is an increasingly important problem Becoming more relevant in the big data era

Proposed a framework for ER in big data Systematic study and organization of all proposed techniques effectiveness layer, meta-blocking, efficiency layer Novel techniques targeted to big data scenarios Lead (up) to more than 95% RR for more than 99% PC!

Several challenges ahead Scalability is always a goal Incremental methods are necessary



collaborations!





Trento, Italy



dbTrento



dbTrento: http://db.disi.unitn.eu



dbTrento



dbTrento

activities described in invited paper: SIGMOD Record, Sep 2012







dbTrento: http://db.disi.unitn.eu





thank you!





[Aizawa et. al., WIRI 2005] Akiko N. Aizawa, Keizo Oyama, "A Fast Linkage Detection Scheme for Multi-Source Information Integration" in WIRI, 2005.

[Baxter et. al., KDD 2003] R. Baxter, P. Christen, T. Churches, "A comparison of fast blocking methods for record linkage", in Workshop on Data Cleaning, Record Linkage and Object Consolidation at KDD, 2003.

[Bilenko et. al., ICDM 2006] M. Bilenko, B. Kamath, R. Mooney, "Adaptive blocking: Learning to scale up record linkage", in ICDM, 2006.

[Christen, TKDE 2011] P. Christen, " A survey of indexing techniques for scalable record linkage and deduplication." in TKDE 2011.

[de Vries et. al., CIKM 2009] T. de Vries, H. Ke, S. Chawla, P. Christen, "Robust record linkage blocking using sux arrays", in CIKM, 2009.

[Gravano et. al., VLDB 2001] L. Gravano, P. Ipeirotis, H. Jagadish, N. Koudas, S. Muthukrishnan, D. Srivastava, "Approximate string joins in a database (almost) for free', in VLDB, 2001.

[Hernandez et. al., SIGMOD 1995] M. Hernandez, S. Stolfo, "The merge/purge problem for large databases", in SIGMOD, 1995.



[Jin et. al., DASFAA 2003] L. Jin, C. Li, S. Mehrotra, "Efficient record linkage in large data sets", in DASFAA, 2003.

[Kim et. al., EDBT 2010] H. Kim, D. Lee, "HARRA: fast iterative hashed record linkage for large-scale data collections", in EDBT, 2010.

[Madhavan et. al., CIDR 2007] J. Madhavan, S. Cohen, X. Dong, A. Halevy, S. Jeffery, D. Ko, C.Yu, "Web-scale data integration: You can afford to pay as you go", in CIDR, 2007

[McCallum et. al., KDD 2000] A. McCallum, K. Nigam, L. Ungar, "Efficient clustering of high-dimensional data sets with application to reference matching", in KDD, 2000.

[Michelson et. al., AAAI 2006] M. Michelson, C. Knoblock, "Learning blocking schemes for record linkage", in AAAI, 2006.

[Papadakis et al., iiWAS 2010] G. Papadakis, G. Demartini, P. Fankhauser, P. Karger, "The missing links: discovering hidden same-as links among a billion of triples", in iiWAS 2010.

[Papadakis et al., WSDM 2011] G. Papadakis, E. Ioannou, C. Niederee, P. Fankhauser, "Efficient entity resolution for large heterogeneous information spaces", in WSDM 2011.



[Papadakis et al., JCDL 2011] G. Papadakis, E. Ioannou, C. Niederee, T. Palpanas, W. Nejdl, "Eliminating the redundancy in blocking-based entity resolution methods", in JCDL 2011.

[Papadakis et al., SWIM 2011] G. Papadakis, E. Ioannou, C. Niederee, T. Palpanas, W. Nejdl, "To Compare or Not to Compare: making Entity Resolution more Efficient", in SWIM workshop (collocated with SIGMOD), 2011.

[Papadakis et al., WSDM 2012] G. Papadakis, E. Ioannou, C. Niederee, T. Palpanas, W. Nejdl, "Beyond 100 million entities: large-scale blocking-based resolution for heterogeneous data", in WSDM 2012.

[Papadakis et. al., TKDE2013] George Papadakis, Ekaterini Ioannou, Themis Palpanas, Claudia Niederee, Wolfgang Nejdl, "A Blocking Framework for Entity Resolution in Highly Heterogeneous Information Spaces", in IEEE TKDE (to appear).

[Papadakis et. al., TKDE] George Papadakis, Georgia Koutrika, Themis Palpanas, Wolfgang Nejdl, "Meta-Blocking: Taking Entity Resolution to the Next Level", in IEEE TKDE (currently under revision).

[Yan et. Al., JCDL 2007] Su Yan, Dongwon Lee, Min-Yen Kan, C. Lee Giles, "Adaptive sorted neighborhood methods for efficient record linkage", in JCDL 2007.