

Entity Resolution for Big Data

Lise Getoor

*University of Maryland
College Park, MD*

Ashwin Machanavajjhala

*Duke University
Durham, NC*

http://www.cs.umd.edu/~getoor/Tutorials/ER_KDD2013.pdf

<http://goo.gl/7tKiiL>

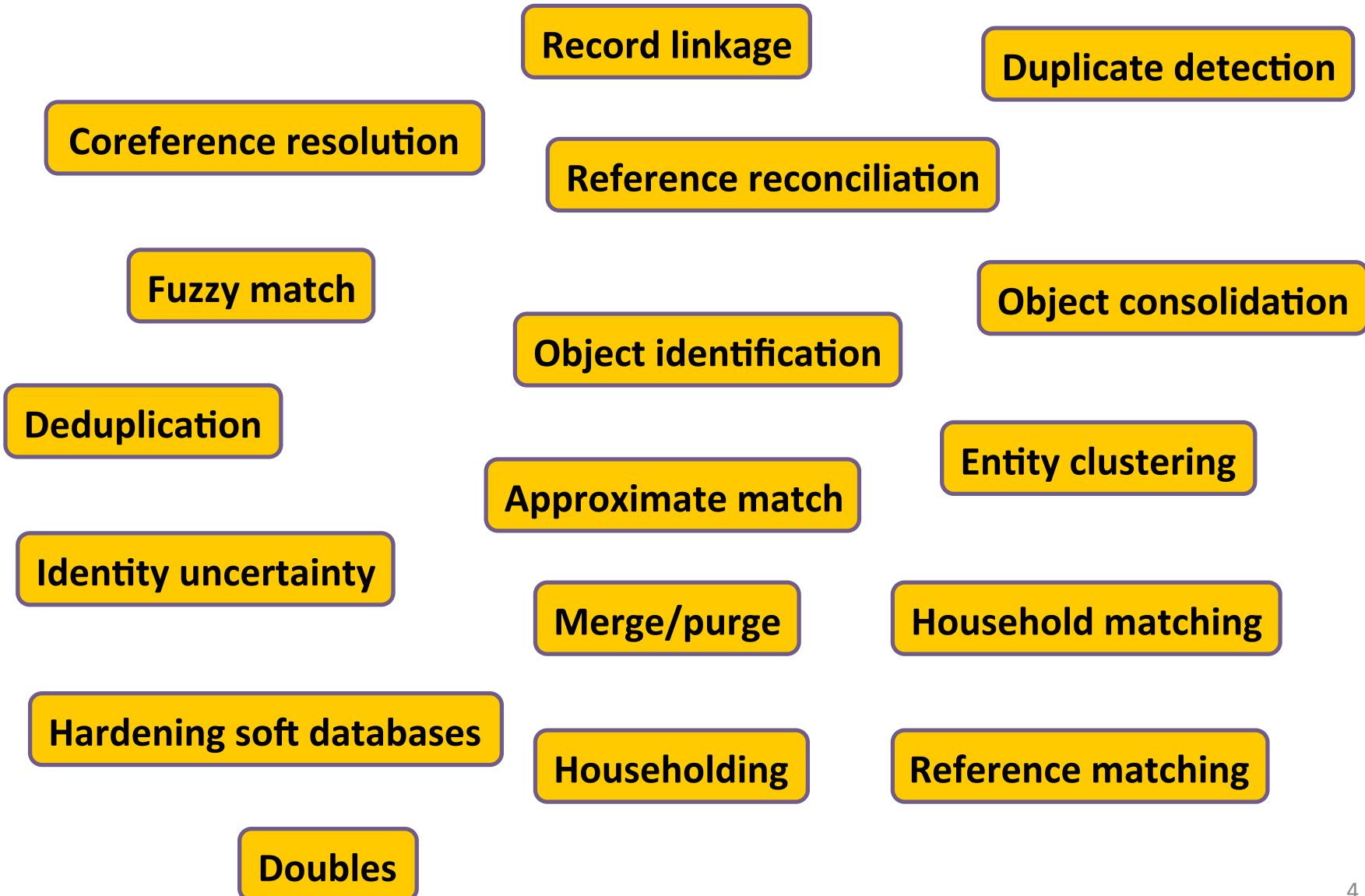
What is Entity Resolution?

Problem of identifying and linking/grouping different manifestations of the same real world object.

Examples of manifestations and objects:

- Different ways of addressing (names, email addresses, FaceBook accounts) the same person in text.
- Different digital photos of the same object.
- Web pages with differing descriptions of the same business.
- ...

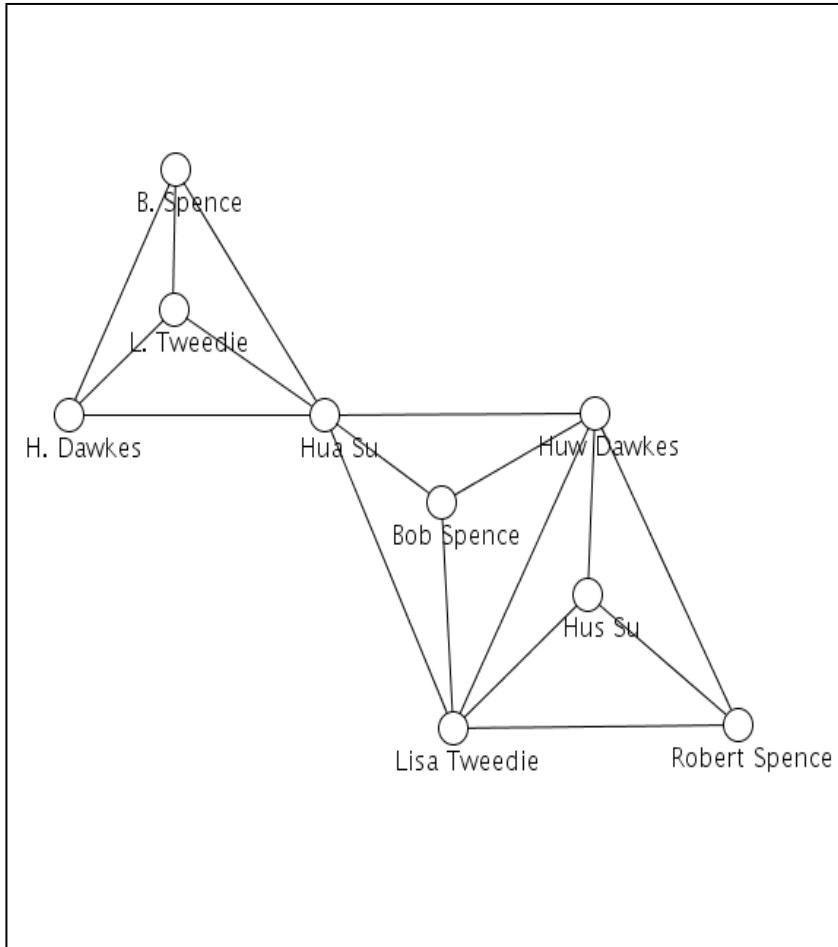
Ironically, Entity Resolution has many duplicate names



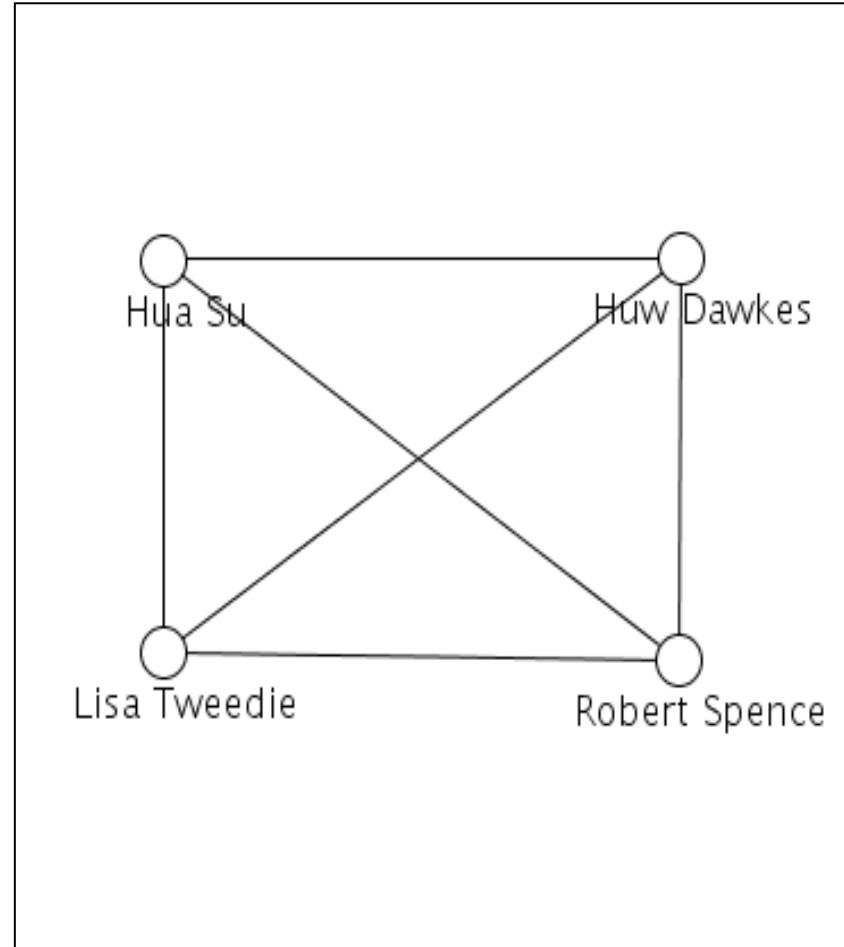
ER Motivating Examples

- *Linking Census Records*
- *Public Health*
- *Medical records*
- *Web search*
- *Comparison shopping*
- *Maintaining customer databases*
- *Law enforcement and Counter-terrorism*
- *Scientific data*
- *Genealogical data*
- *Bibliographic data*

Motivation: ER and Network Analysis



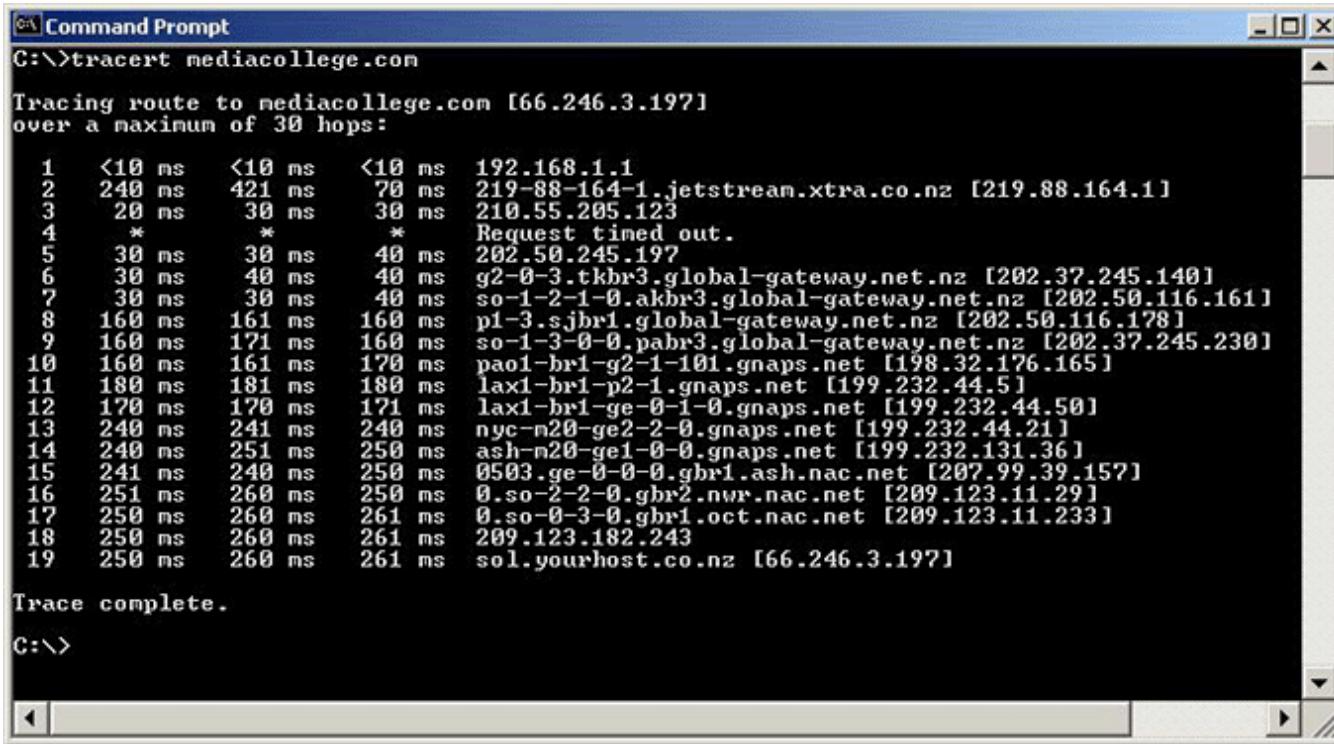
before



after

Motivation: ER and Network Analysis

- Measuring the topology of the internet ... using traceroute



The screenshot shows a Windows Command Prompt window titled "Command Prompt". The command entered is "C:\>tracert mediacollege.com". The output displays the traceroute results to the website mediacollege.com, which has the IP address 66.246.3.197. The route consists of 19 hops, with the last hop being the destination. Hops 4 through 19 show varying levels of network latency, with some segments taking up to 240 ms. The output also includes a note about a request timing out at hop 4.

```
C:\>tracert mediacollege.com
Tracing route to mediacollege.com [66.246.3.197]
over a maximum of 30 hops:
1 <10 ms <10 ms <10 ms 192.168.1.1
2 240 ms 421 ms 70 ms 219-88-164-1.jetstream.xtra.co.nz [219.88.164.1]
3 20 ms 30 ms 30 ms 210.55.205.123
4 * * * Request timed out.
5 30 ms 30 ms 40 ms 202.50.245.197
6 30 ms 40 ms 40 ms g2-0-3.tkbr3.global-gateway.net.nz [202.37.245.140]
7 30 ms 30 ms 40 ms so-1-2-1-0.akbr3.global-gateway.net.nz [202.50.116.161]
8 160 ms 161 ms 160 ms p1-3.sjbr1.global-gateway.net.nz [202.50.116.178]
9 160 ms 171 ms 160 ms so-1-3-0-0.pabr3.global-gateway.net.nz [202.37.245.230]
10 160 ms 161 ms 170 ms pa01-br1-g2-1-101.gnaps.net [198.32.176.165]
11 180 ms 181 ms 180 ms lax1-br1-p2-1.gnaps.net [199.232.44.5]
12 170 ms 170 ms 171 ms lax1-br1-ge-0-1-0.gnaps.net [199.232.44.50]
13 240 ms 241 ms 240 ms nyc-n20-ge2-2-0.gnaps.net [199.232.44.21]
14 240 ms 251 ms 250 ms ash-n20-ge1-0-0.gnaps.net [199.232.131.36]
15 241 ms 240 ms 250 ms 0503.ge-0-0-0.gbr1.ash.nac.net [207.99.39.157]
16 251 ms 260 ms 250 ms 0.so-2-2-0.gbr2.nwr.nac.net [209.123.11.29]
17 250 ms 260 ms 261 ms 0.so-0-3-0.gbr1.oct.nac.net [209.123.11.233]
18 250 ms 260 ms 261 ms 209.123.182.243
19 250 ms 260 ms 261 ms sol.yourhost.co.nz [66.246.3.197]

Trace complete.
C:\>
```

IP Aliasing Problem [Willinger et al. 2009]

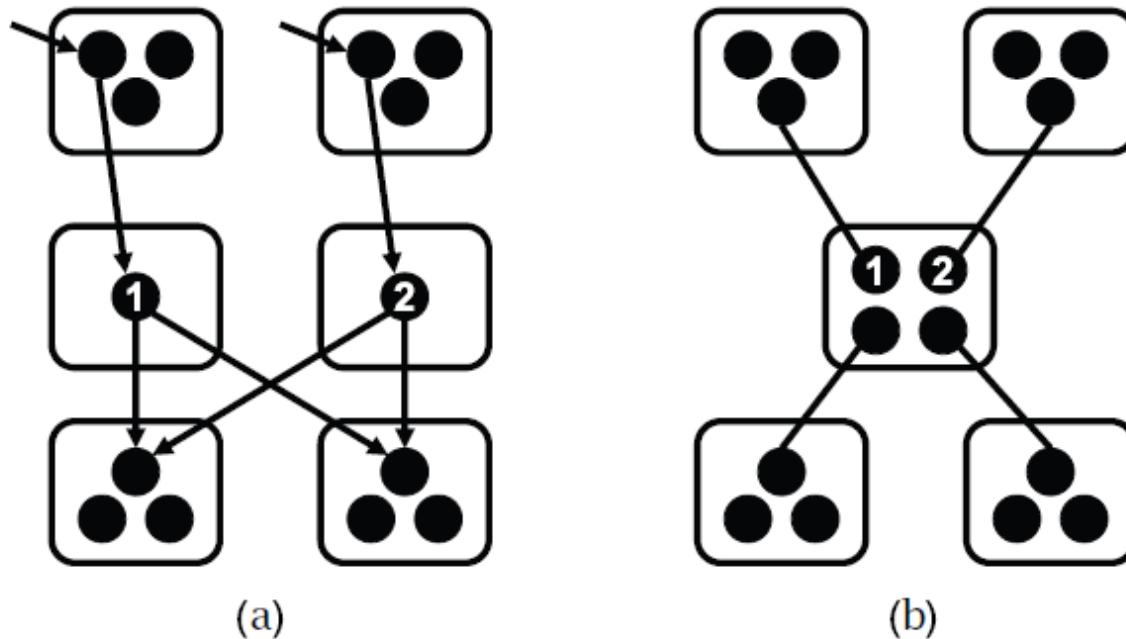


Figure 2. The IP alias resolution problem.
Paraphrasing Fig. 4 of [50], traceroute does not list routers (boxes) along paths but IP addresses of input interfaces (circles), and alias resolution refers to the correct mapping of interfaces to routers to reveal the actual topology. In the case where interfaces 1 and 2 are aliases, (b) depicts the actual topology while (a) yields an “inflated” topology with more routers and links than the real one.

IP Aliasing Problem [Willinger et al. 2009]

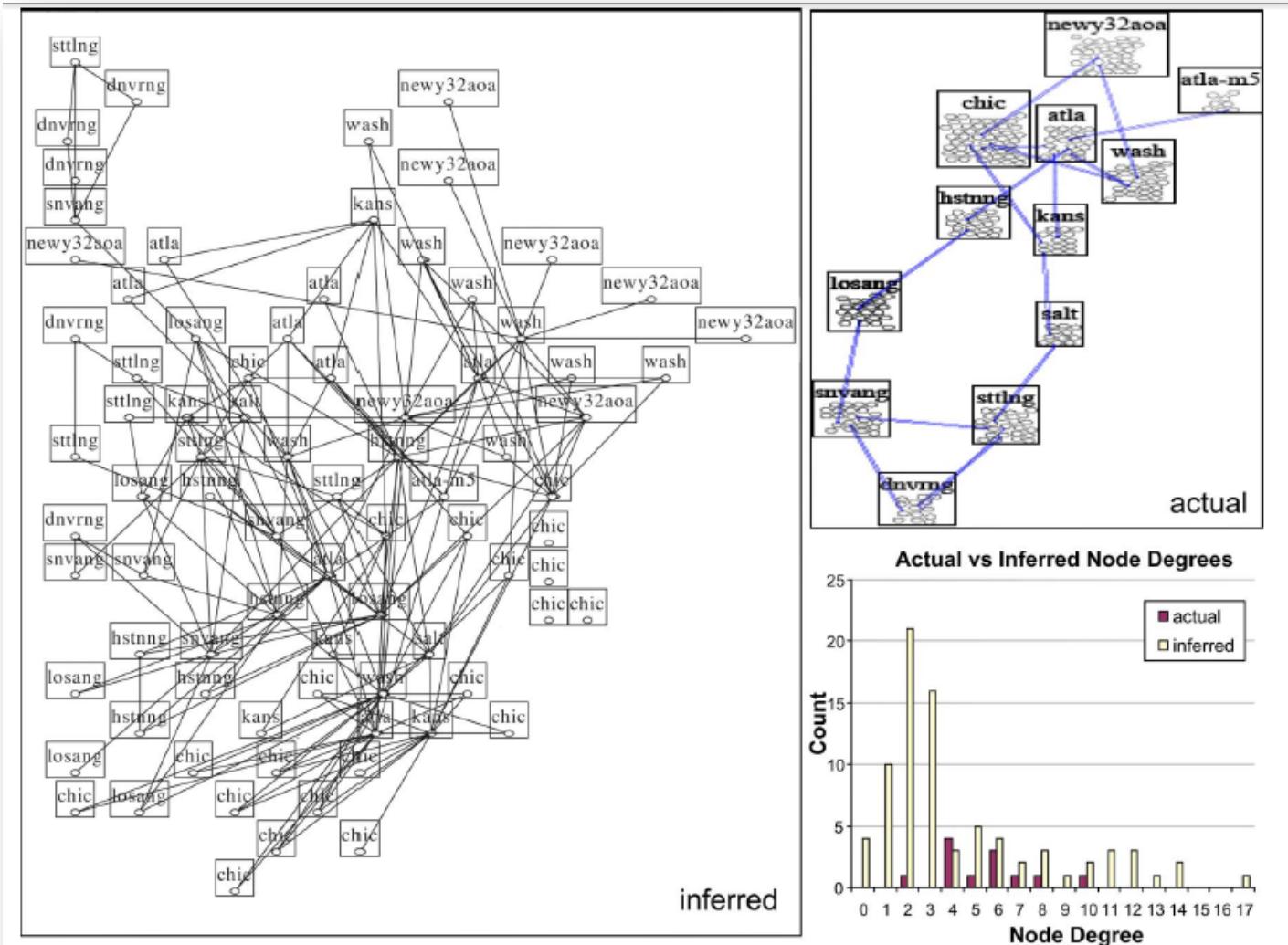
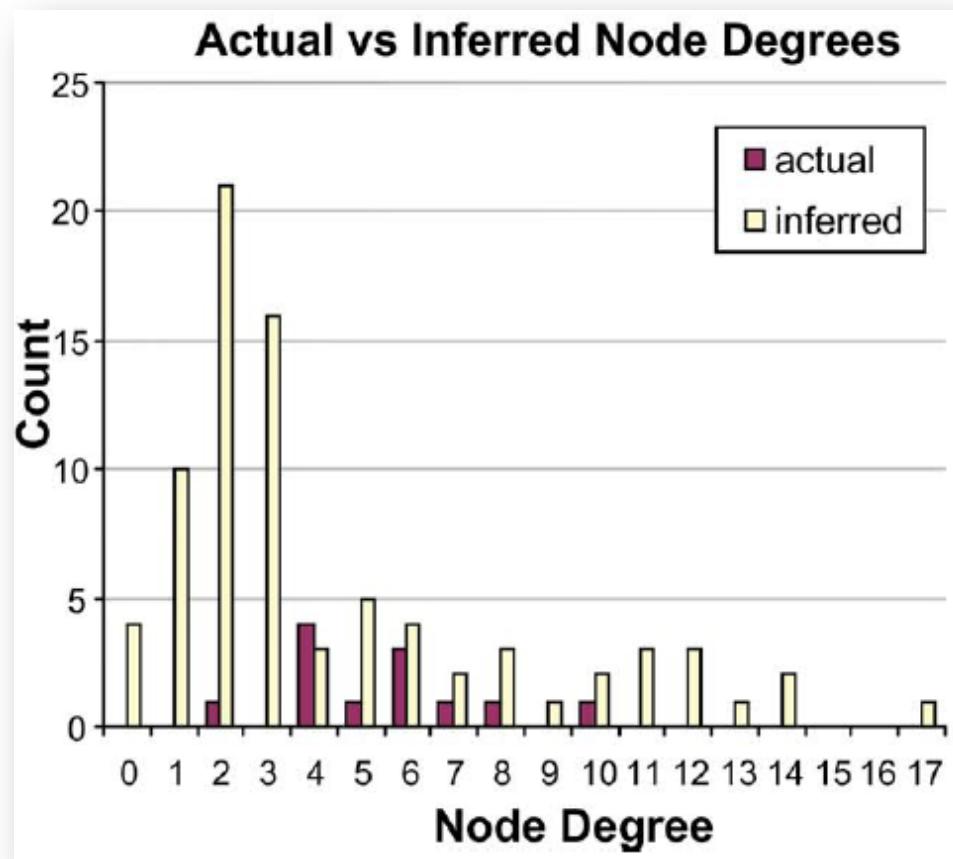


Figure 3. The IP alias resolution problem in practice. This is re-produced from [48] and shows a comparison between the Abilene/Internet2 topology inferred by Rocketfuel (left) and the actual topology (top right). Rectangles represent routers with interior ovals denoting interfaces. The histograms of the corresponding node degrees are shown in the bottom right plot. © 2008 ACM,

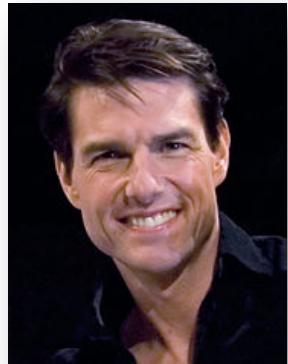
IP Aliasing Problem [Willinger et al. 2009]



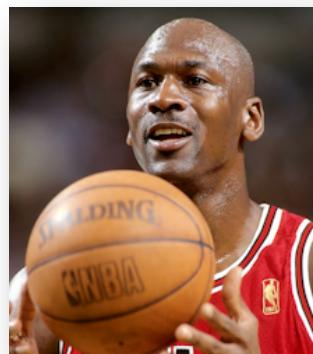
Traditional Challenges in ER

- Name/Attribute ambiguity

Thomas Cruise



Michael Jordan



Traditional Challenges in ER

- Name/Attribute ambiguity
- Errors due to data entry



	C1	C2
	Total Cholesterol_1	Total Cholesterol_2
682	214.4	214.4
683	184.4	184.4
684	183.5	183.5
685	240.7	240.7
686	215.1	215.1
687	198.6	198.6
688	2800.0	280.0
689	210.8	210.8
690	182.5	182.5
691	192.6	192.6

Traditional Challenges in ER

- Name/Attribute ambiguity
- Errors due to data entry
- Missing Values

Exhibit 2: Examples of variables that are set to unknown values

Administrative dates: set to 0101YY, 010199, 999999

Date of Birth 0101YY, 1506YY, 3006YY, 0107YY, 1507YY, 0101YEAR

Names: set to spaces, NK, UNKNOWN, or ZZZZ
BABY, MALE, FEMALE, TWIN, TRIPLET, INFANT

Other variables: set to 9, 99, 9999, -1
NK (Not Known)
NA (Not applicable)
NC (Not coded)
U (Unknown)

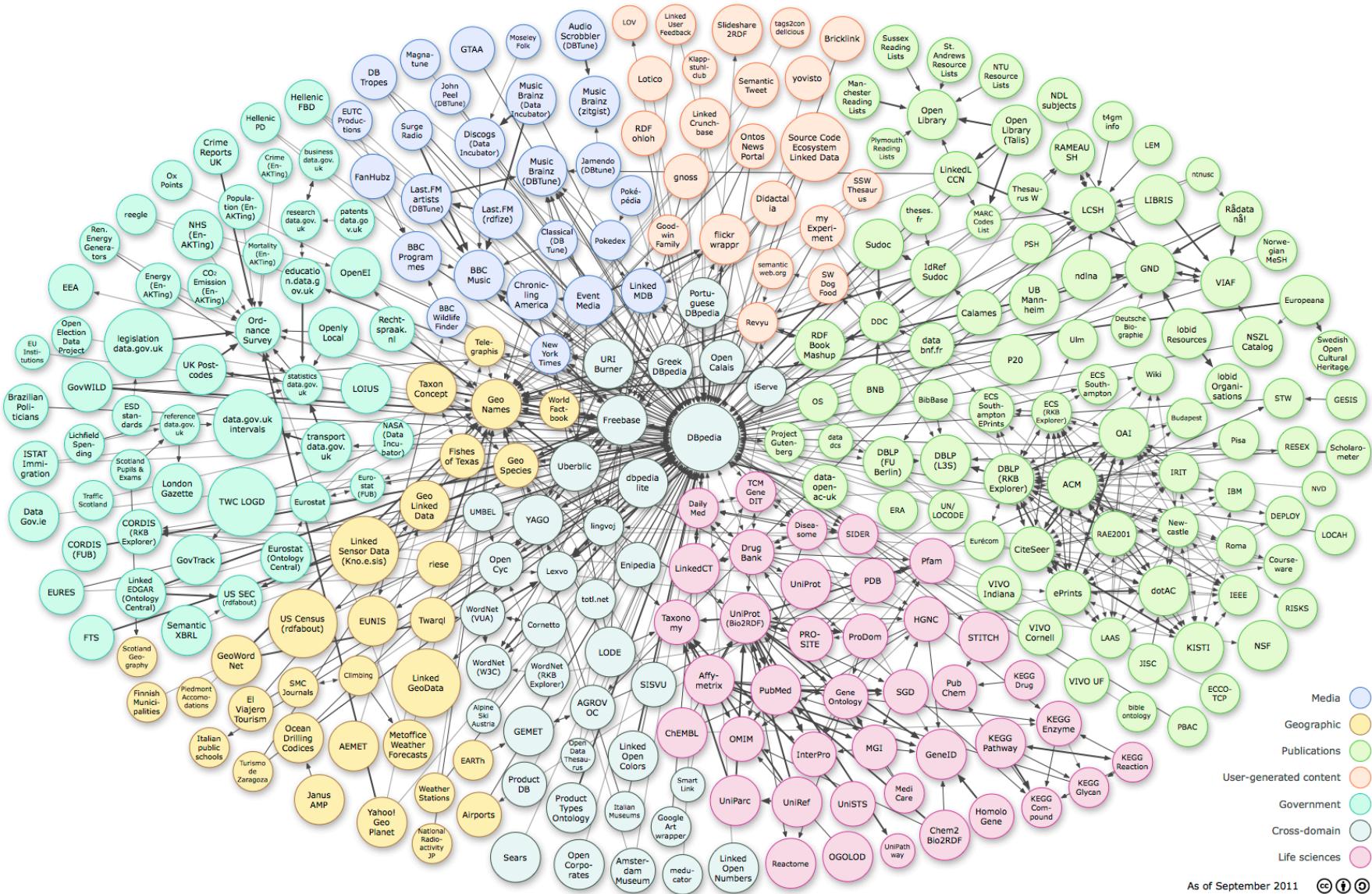
Traditional Challenges in ER

- Name/Attribute ambiguity
- Errors due to data entry
- Missing Values
- Changing Attributes



- Data formatting
- Abbreviations / Data Truncation

Big-Data ER Challenges



Big-Data ER Challenges

- Larger and more Datasets
 - Need efficient parallel techniques
- More Heterogeneity
 - Unstructured, Unclean and Incomplete data. Diverse data types.
 - No longer just matching names with names, but Amazon profiles with browsing history on Google and friends network in Facebook.

Big-Data ER Challenges

- Larger and more Datasets
 - Need efficient parallel techniques
- More Heterogeneity
 - Unstructured, Unclean and Incomplete data. Diverse data types.
- More linked
 - Need to infer relationships in addition to “equality”
- Multi-Relational
 - Deal with structure of entities (Are Walmart and Walmart Pharmacy the same?)
- Multi-domain
 - Customizable methods that span across domains
- Multiple applications (web search versus comparison shopping)
 - Serve diverse application with different accuracy requirements

Outline

1. Abstract Problem Statement
 - a) Various formulations of the problem
 - b) Typical constraints that appear in ER formulations
 - c) Differences from standard ML tasks
2. Algorithmic Foundations of ER
3. Scaling ER to Big-Data
4. Challenges & Future Directions

Outline

1. Abstract Problem Statement
2. Algorithmic Foundations of ER
 - a) Data Preparation and Match Features
 - b) Pairwise ER
 - c) Algorithms for Enforcing Constraints
 - Record Linkage
 - Deduplication
 - Collective ER
3. Scaling ER to Big-Data
4. Challenges & Future Directions

Break

Outline

1. Abstract Problem Statement
2. Algorithmic Foundations of ER
3. Scaling ER to Big-Data
 - a) Blocking/Canopy Generation
 - b) Distributed ER
4. Challenges & Future Directions

Outline

1. Abstract Problem Statement
2. Algorithmic Foundations of ER
3. Scaling ER to Big-Data
4. Challenges & Future Directions

Scope of the Tutorial

- What we cover:
 - Fundamental algorithmic concepts in ER
 - Scaling ER to big datasets
 - Taxonomy of current ER algorithms
- What we do not cover:
 - Schema/ontology resolution
 - Data fusion/integration/exchange/cleaning
 - Entity/Information Extraction
 - Privacy aspects of Entity Resolution
 - Details on similarity measures
 - Technical details and proofs

ER References

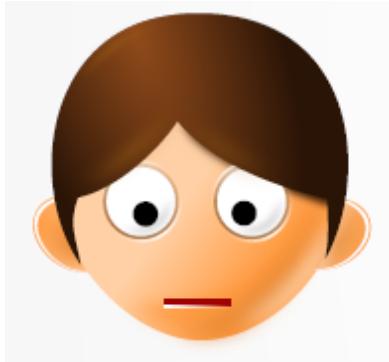
- Book / Survey Articles
 - Data Quality and Record Linkage Techniques
[T. Herzog, F. Scheuren, W. Winkler, Springer, '07]
 - Duplicate Record Detection [A. Elmagrid, P. Ipeirotis, V. Verykios, TKDE '07]
 - An Introduction to Duplicate Detection [F. Naumann, M. Herschel, M&P synthesis lectures 2010]
 - Evaluation of Entity Resolution Approached on Real-world Match Problems
[H. Kopke, A. Thor, E. Rahm, PVLDB 2010]
 - Data Matching [P. Christen, Springer 2012]
- Tutorials
 - Record Linkage: Similarity measures and Algorithms
[N. Koudas, S. Sarawagi, D. Srivatsava SIGMOD '06]
 - Data fusion--Resolving data conflicts for integration
[X. Dong, F. Naumann VLDB '09]
 - Entity Resolution: Theory, Practice and Open Challenges
<http://goo.gl/Ui38o> [L. Getoor, A. Machanavajjhala AAAI '12]

PART 1

ABSTRACT PROBLEM STATEMENT

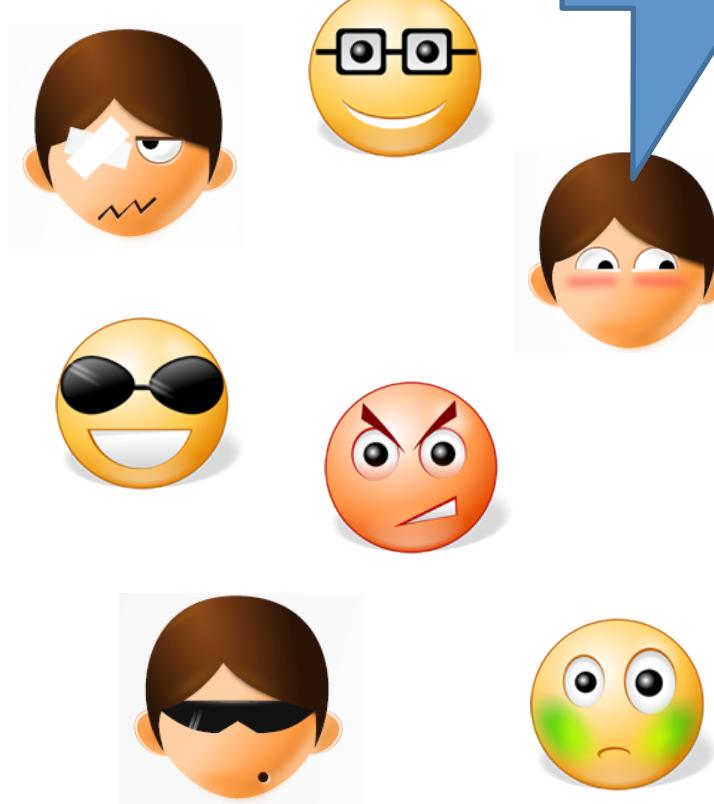
Abstract Problem Statement

Real World



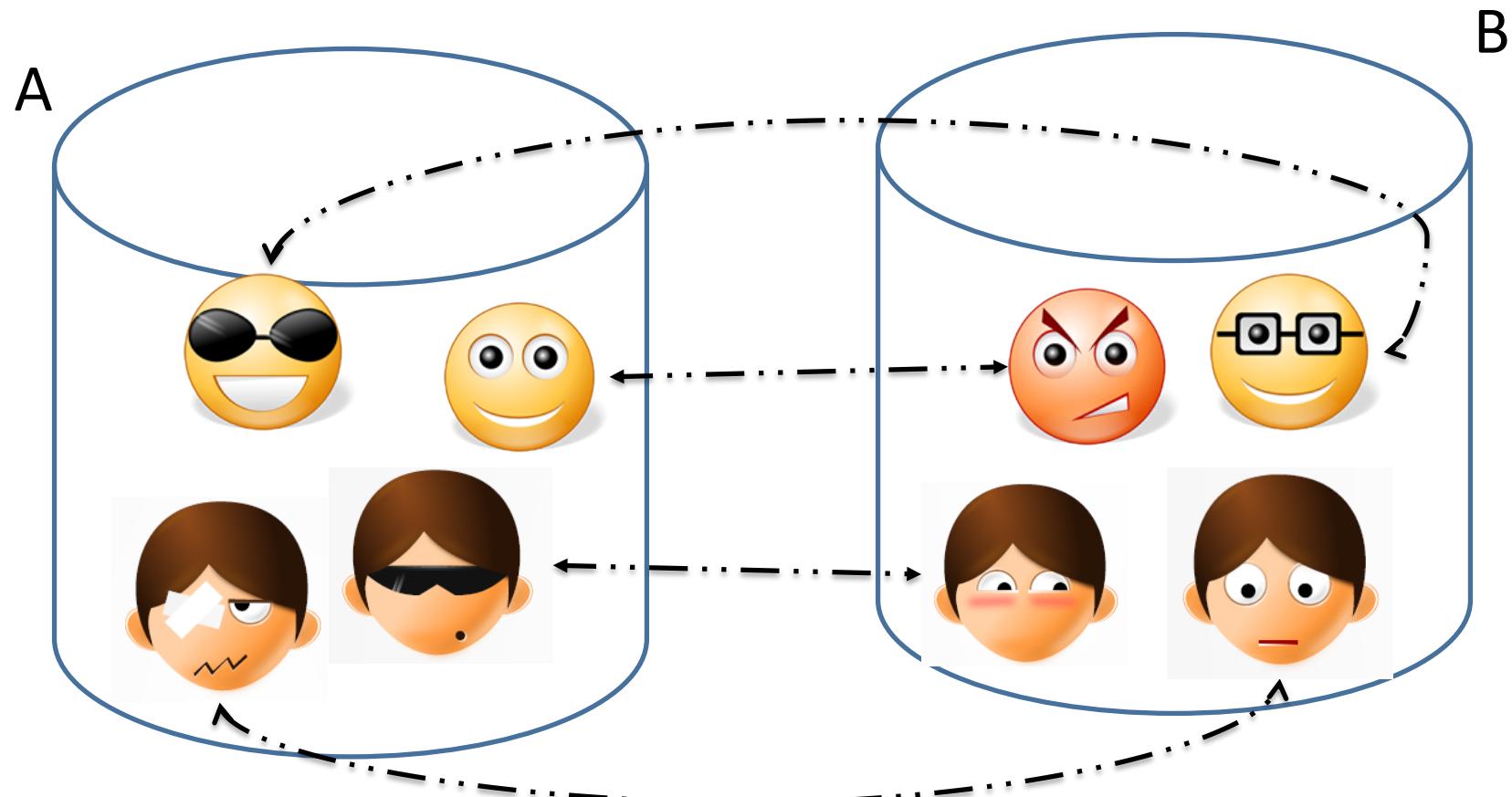
Digital World

Records /
Mentions



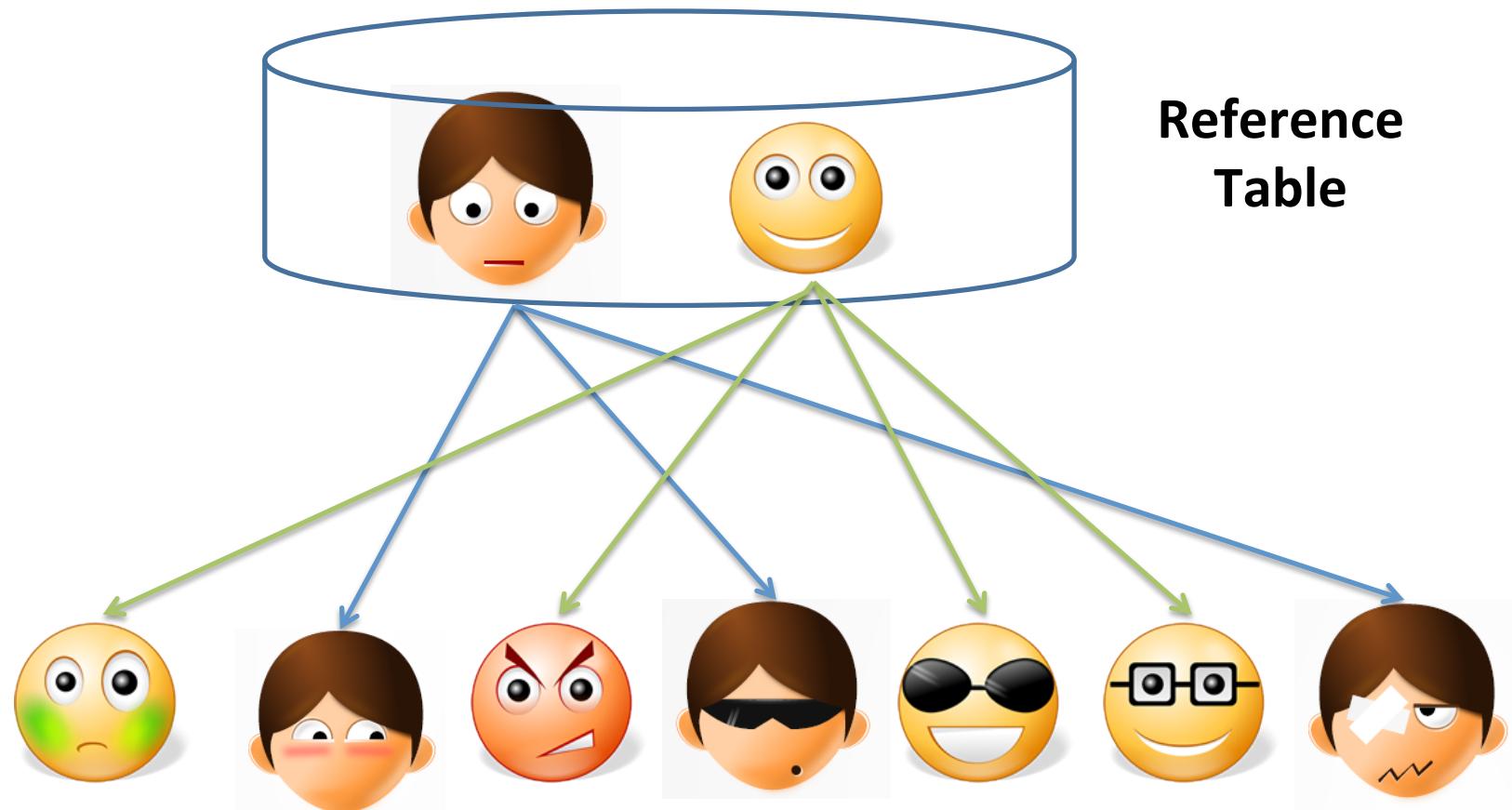
Record Linkage Problem Statement

- Link records that match across databases



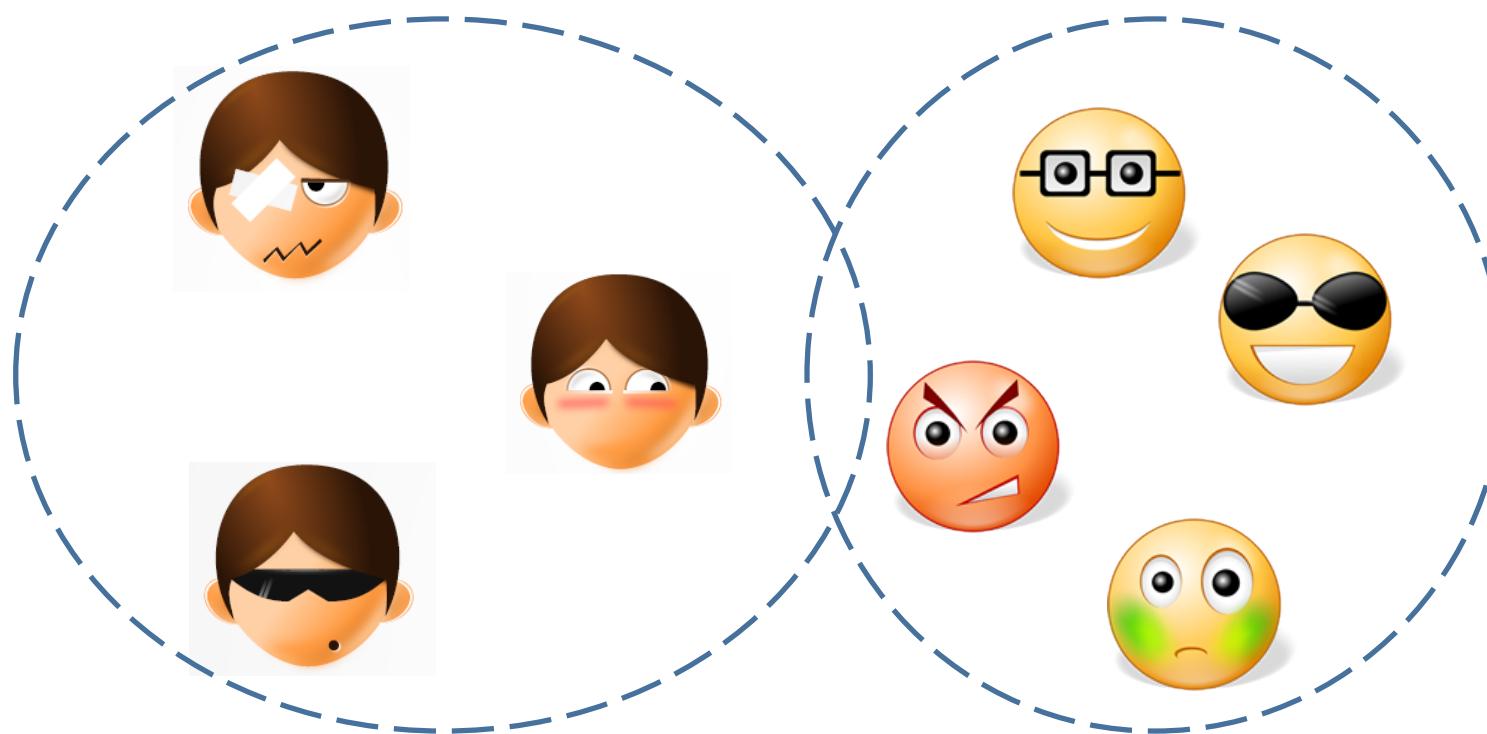
Reference Matching Problem

- Match noisy records to clean records in a reference table



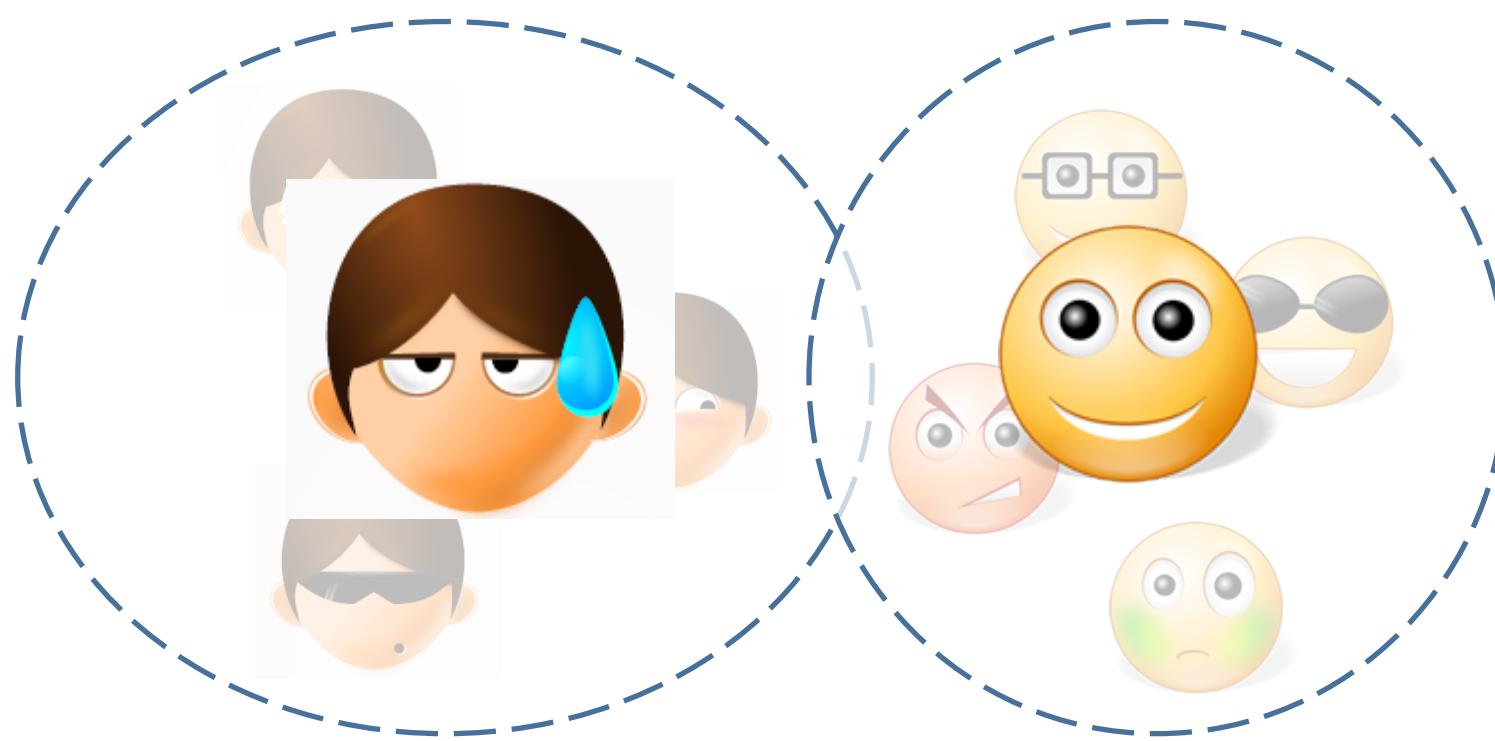
Deduplication Problem Statement

- Cluster the records/mentions that correspond to same entity



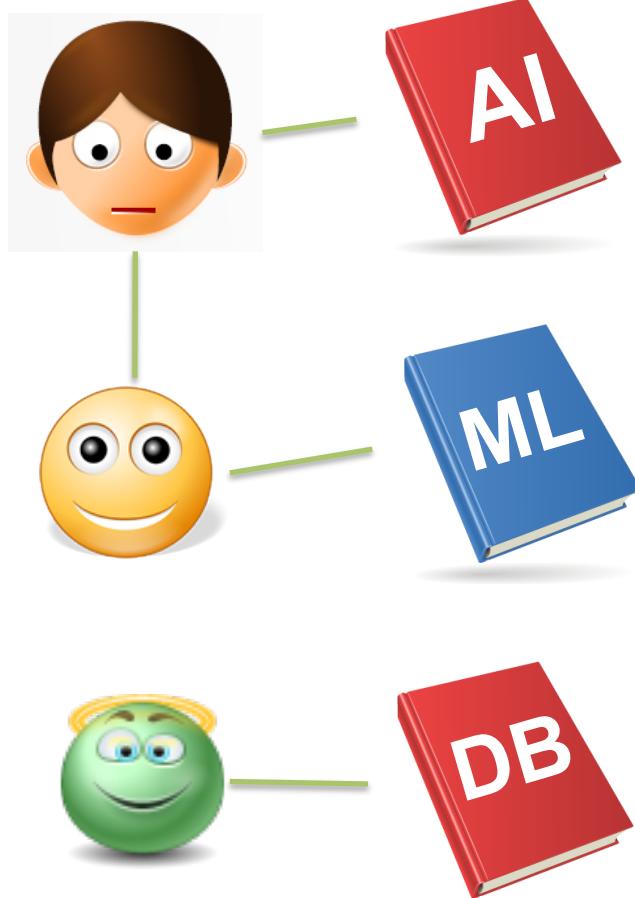
Deduplication Problem Statement

- Cluster the records/mentions that correspond to same entity
 - **Intensional Variant:** Compute cluster representative

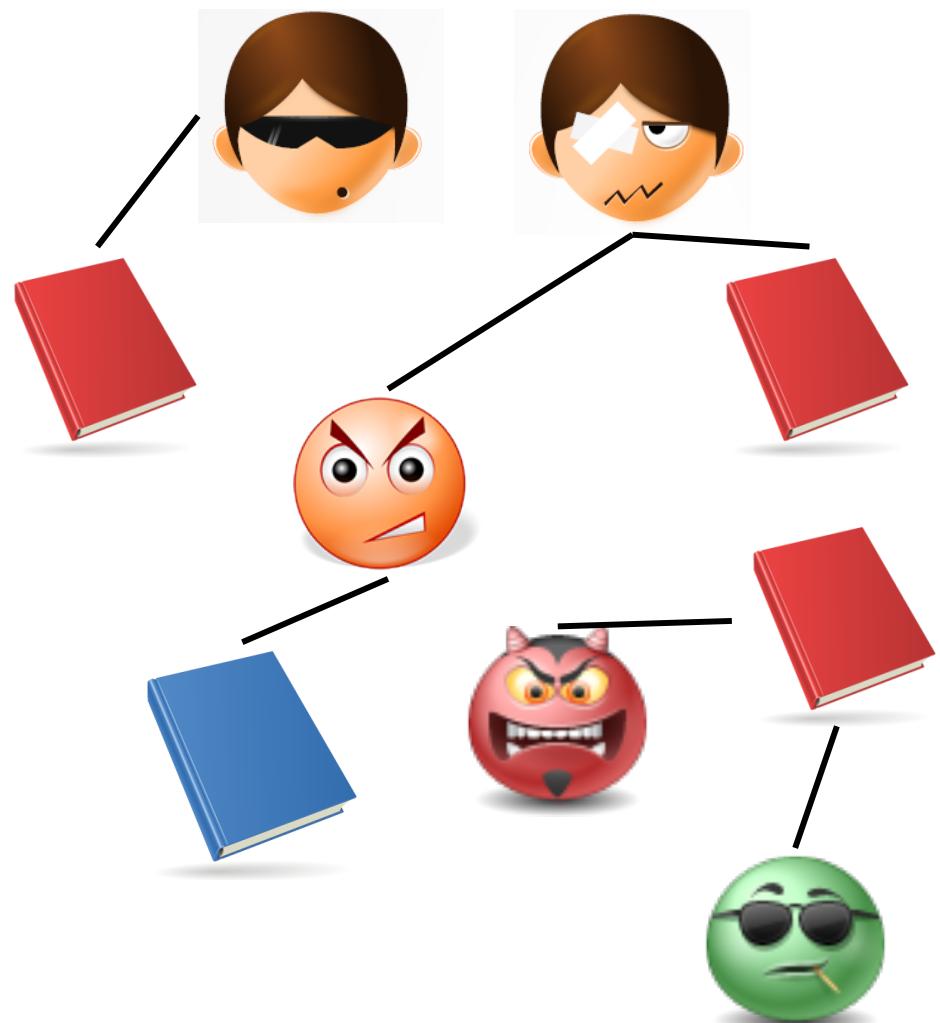


Abstract Problem Statement

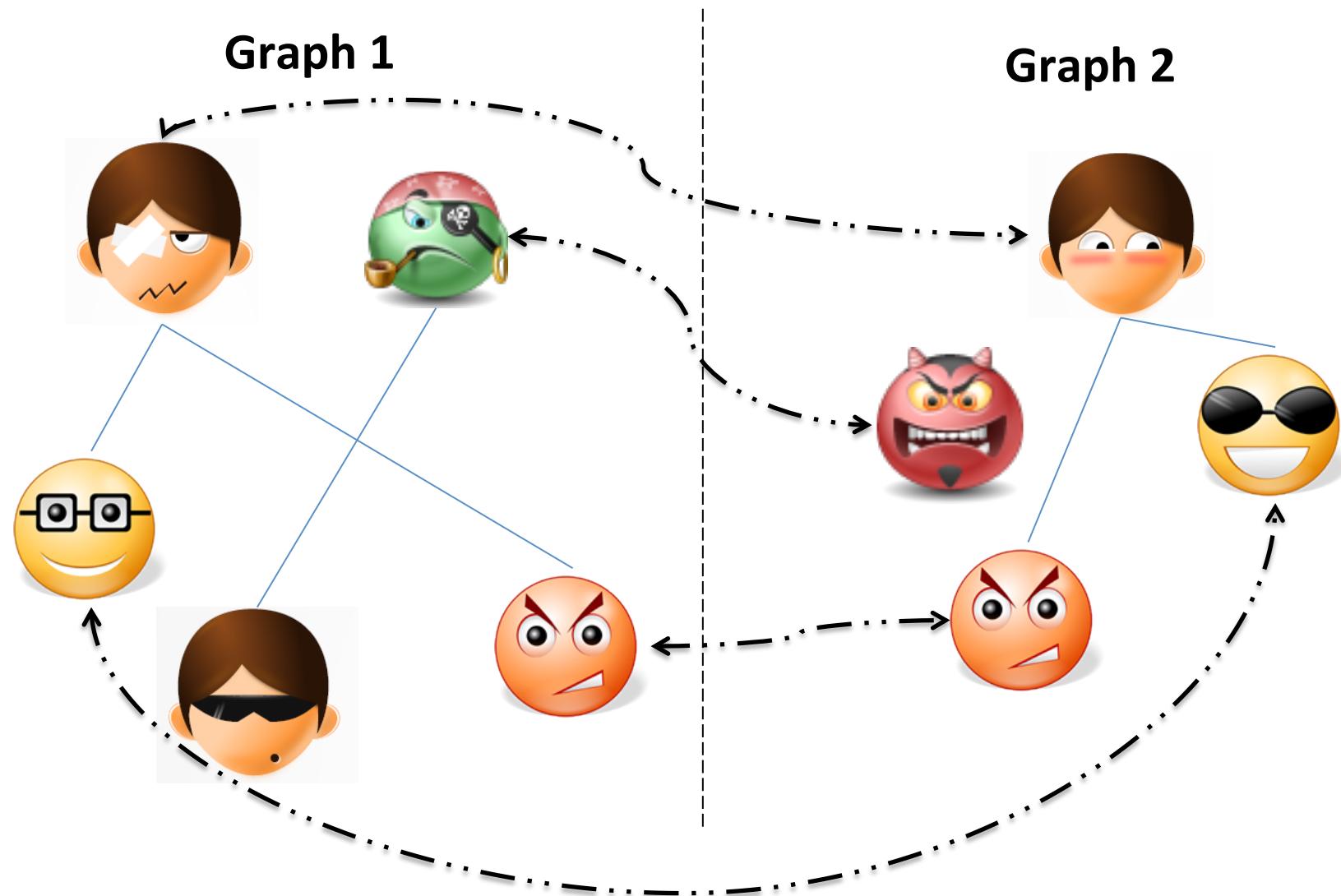
Real World



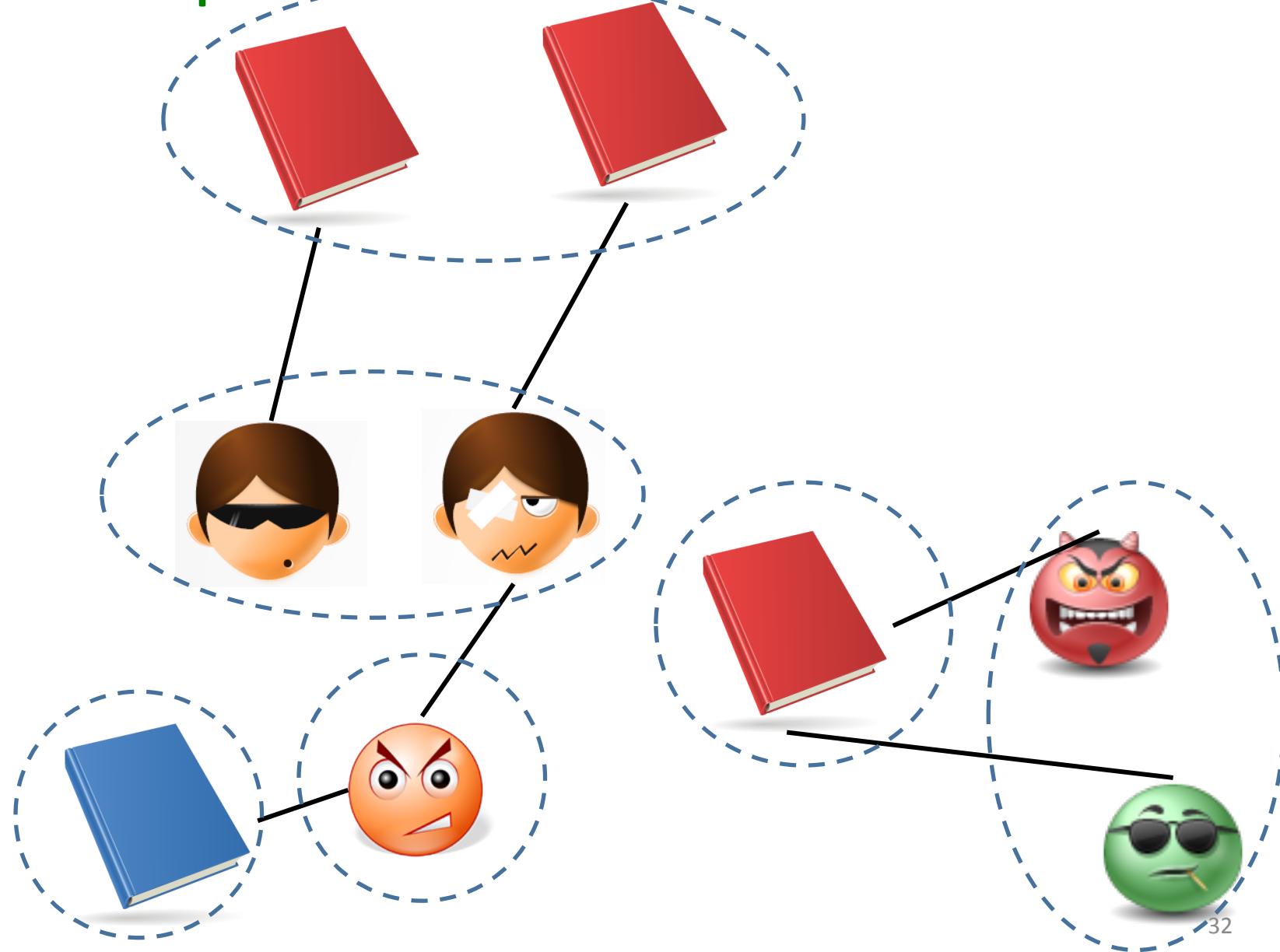
Digital World



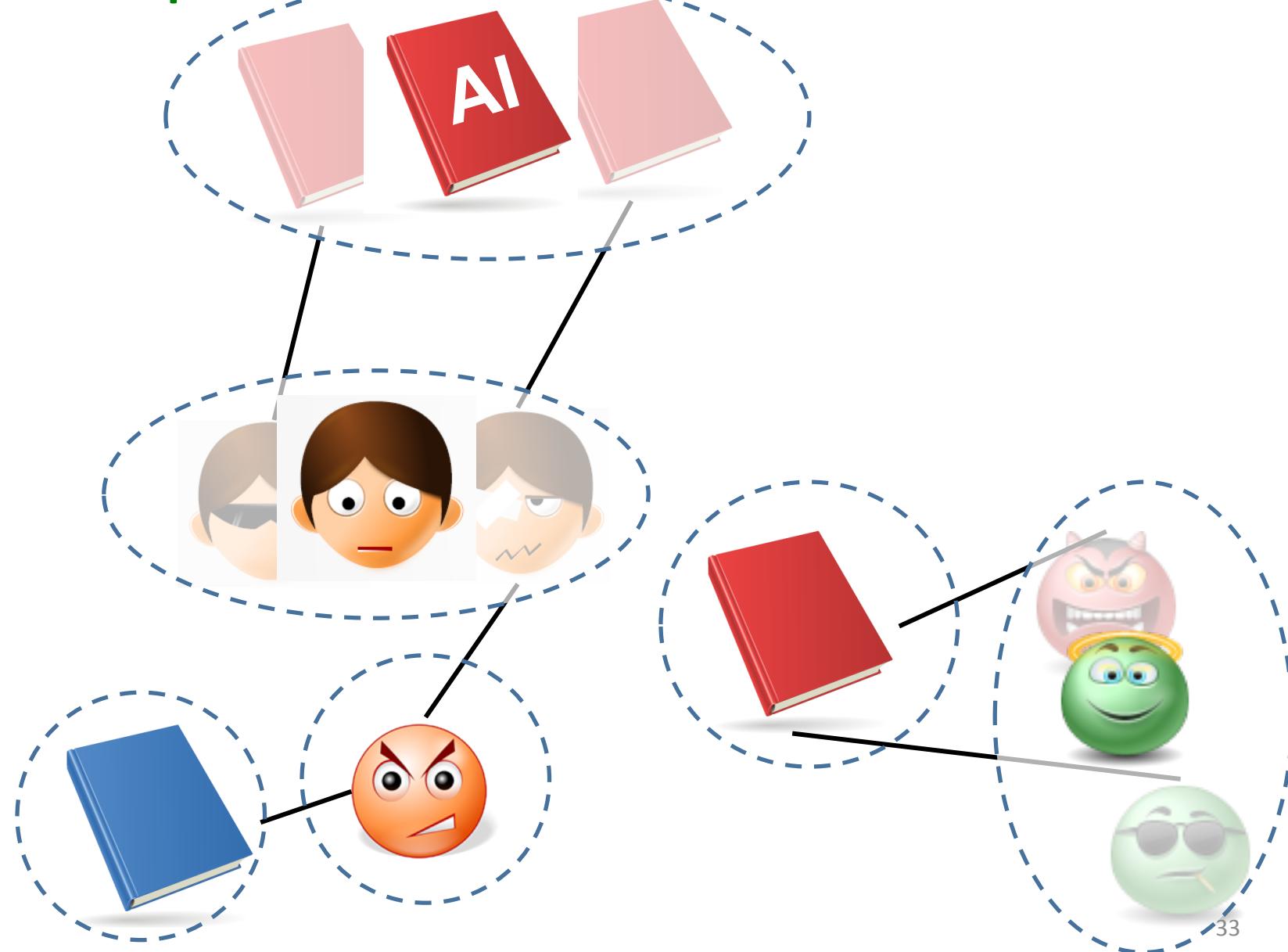
Graph/Motif Alignment



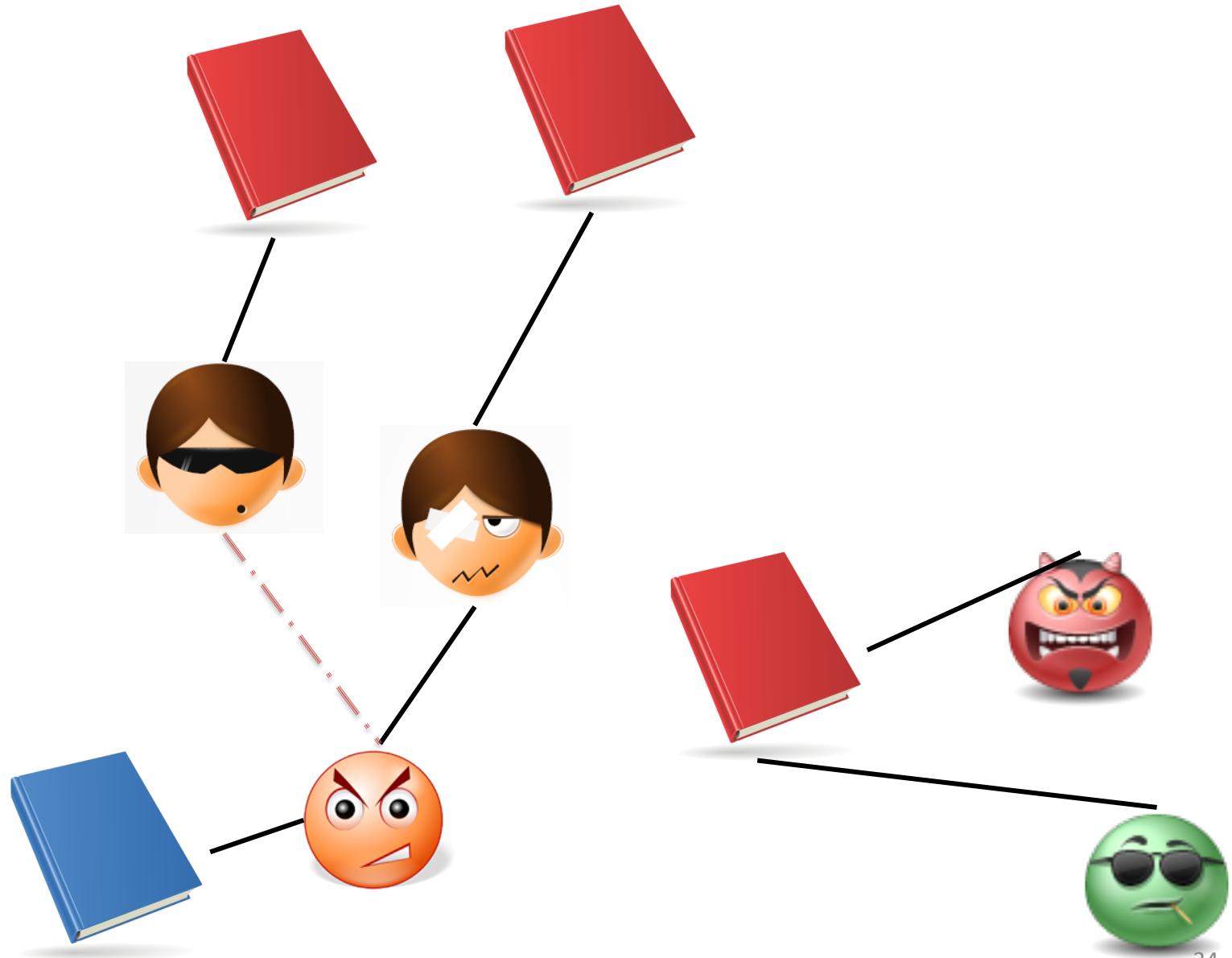
Deduplication Problem Statement



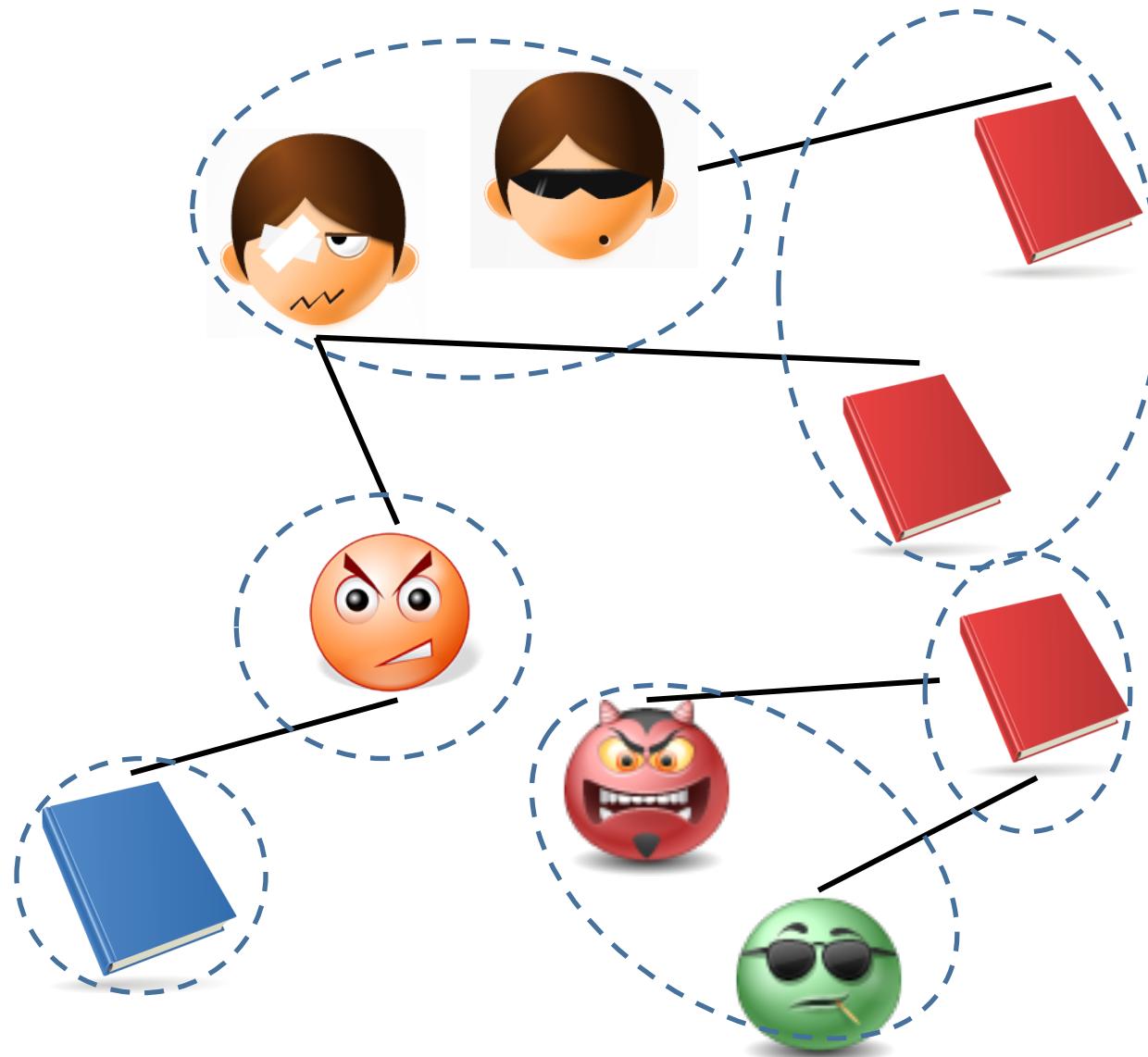
Deduplication with Canonicalization



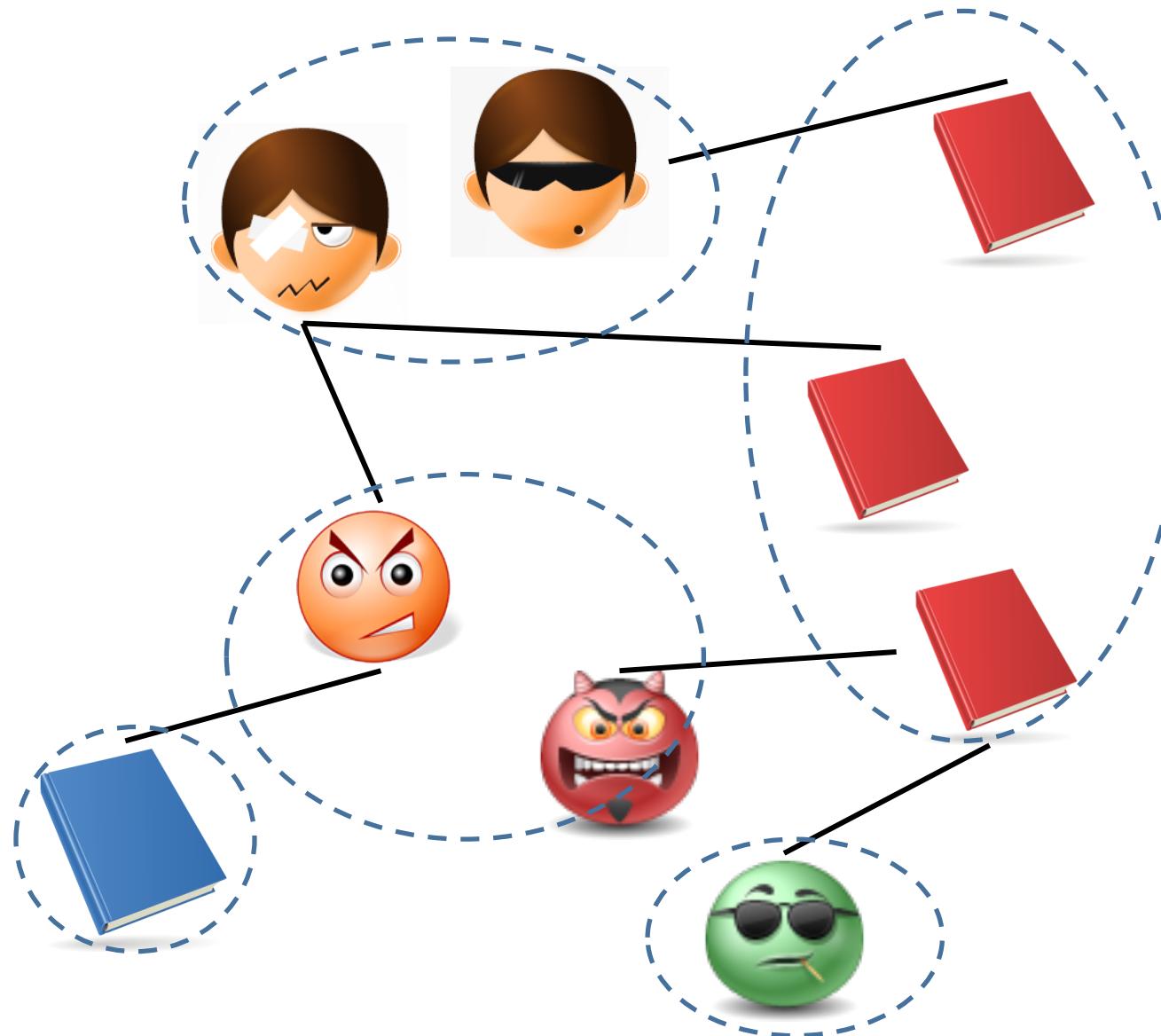
Link Prediction Problem Statement



Relationships are crucial

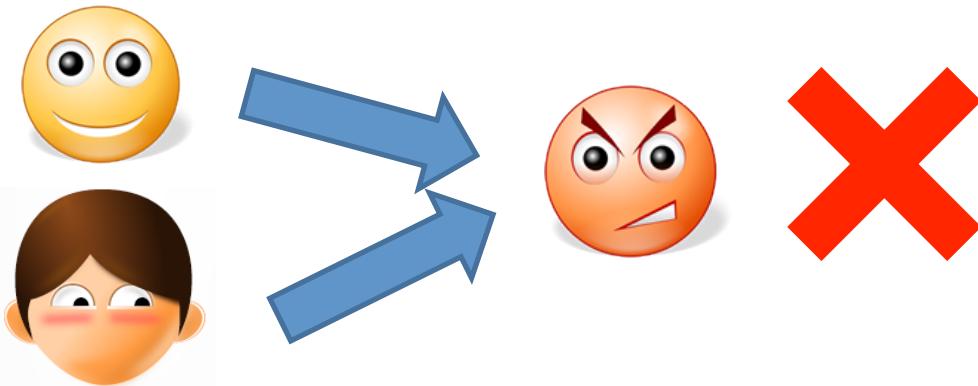


Relationships are crucial



Typical Assumptions Made

- *Each record/mention is associated with a single real world entity.*



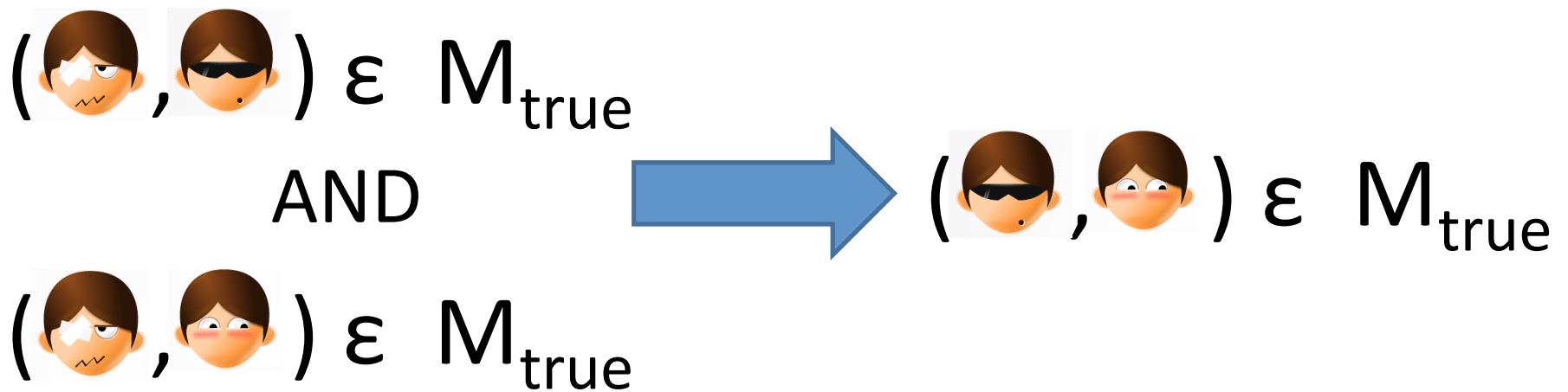
- *If two records/mentions are identical, then they are true matches*

$$(\text{boy with sunglasses}, \text{boy with sunglasses}) = M_{\text{true}}$$

ER versus Classification

Finding matches vs non-matches is a classification problem

- Imbalanced: typically $O(R)$ matches, $O(R^2)$ non-matches
- Instances are pairs of records. Pairs are not IID



ER vs (Multi-relational) Clustering

Computing entities from records is a clustering problem

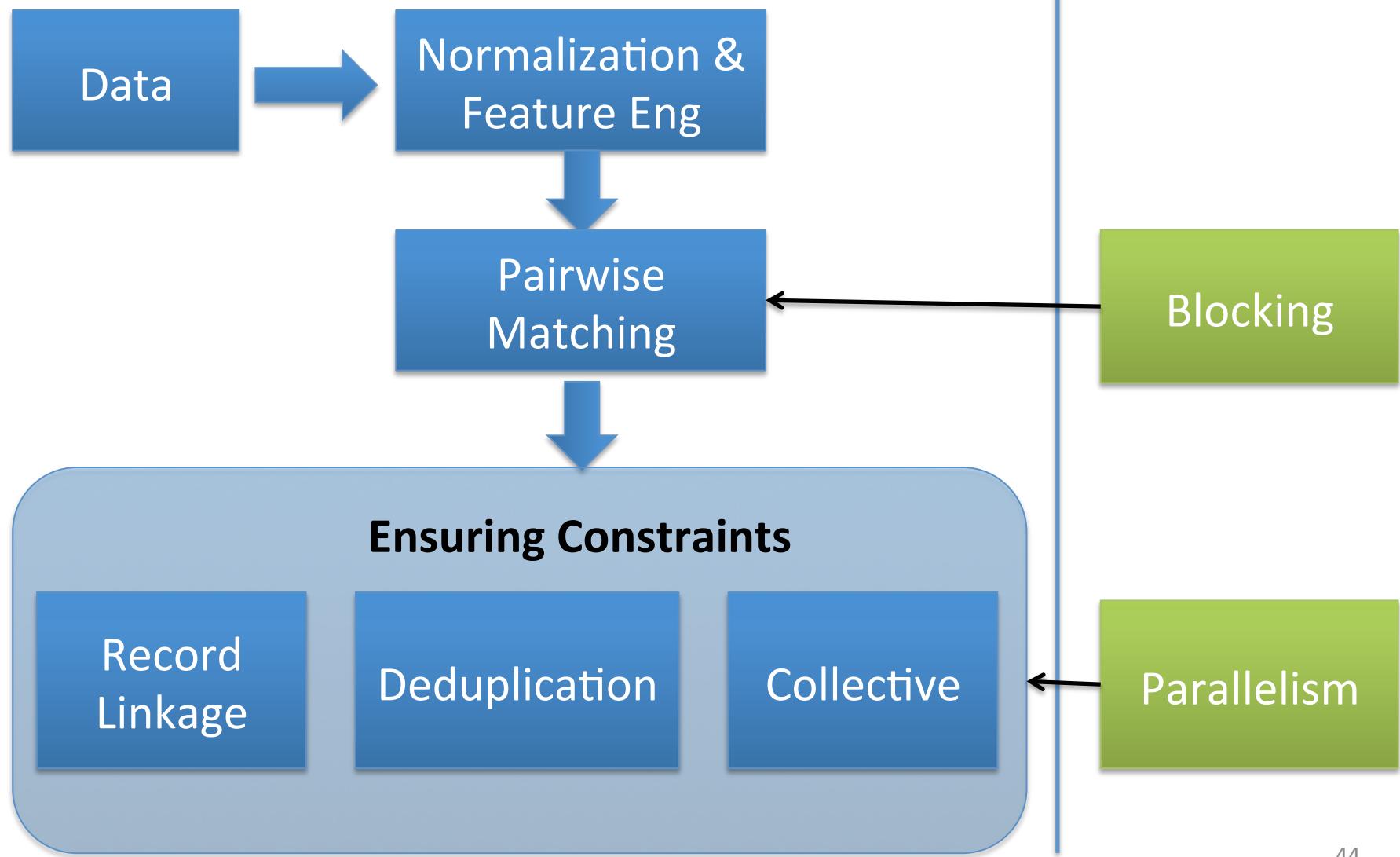
- In typical clustering algorithms (k-means, LDA, etc.) *number of clusters is a constant or sub linear in R.*
- In ER: *number of clusters is linear in R, and average cluster size is a constant. Significant fraction of clusters are singletons.*

Constraints

- Important forms of constraints:
 - **Exclusivity:** If M1 matches with M2, then M3 cannot match with M2
 - **Transitivity:** If M1 and M2 match, M2 and M3 match, then M1 and M3 match
 - **Functional Dependency:** If M1 and M2 match, then M3 and M4 must match
- Exclusivity is key to record linkage
- Transitivity is key to deduplication
- Functional dependencies appear in multi-relational ER

Part 2: Algorithmic Foundations

Part 3: Scaling

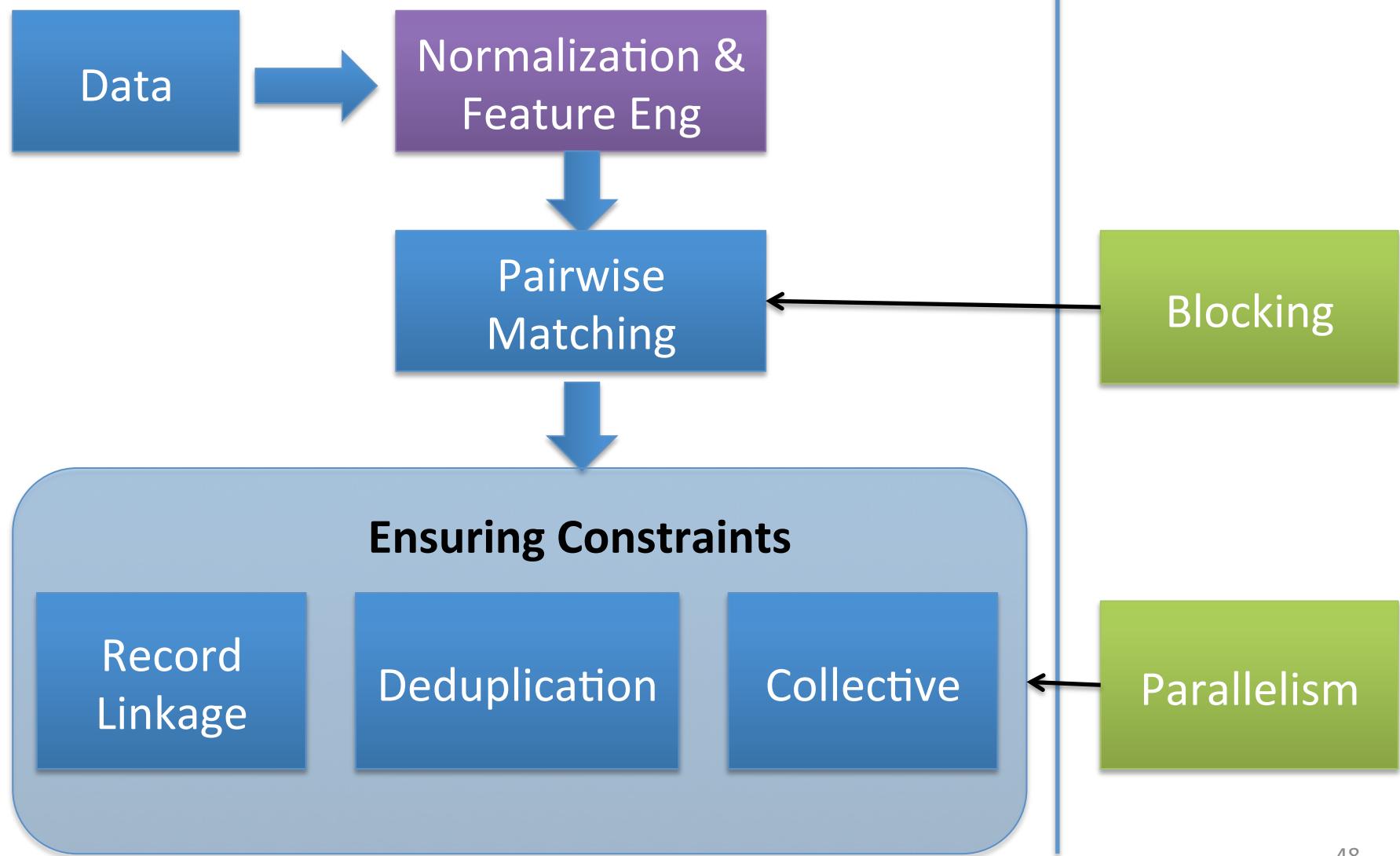


PART 2

ALGORITHMIC FOUNDATIONS OF ER

Part 2: Algorithmic Foundations

Part 3: Scaling



Data Preparation

- Need to perform a number of tasks before ER.
- Entity Identification/Information Extraction
 - Text/semi structured data, Social media, ...
 - Named-entity recognition, segmentation/labeling, source discovery ...
- Feature Engineering
 - Identifying attributes that might be useful for matching
- Schema Normalization
- Data Normalization
- Not a focus of this tutorial

Normalization

- **Schema normalization**
 - Schema Matching – e.g., contact number and phone number
 - Compound attributes – full address vs str, city, state, zip
 - Nested attributes
 - List of features in one dataset (air conditioning, parking) vs each feature a boolean attribute
 - Set valued attributes
 - Set of phones vs primary/secondary phone
 - Record segmentation from text
- **Data normalization**
 - Often convert to all lower/all upper; remove whitespace
 - detecting and correcting values that contain known typographical errors or variations,
 - expanding abbreviations and replacing them with standard forms; replacing nicknames with their proper name forms
 - Usually done based on dictionaries (e.g., commercial dictionaries, postal addresses, etc.)

Normalization

- **Schema normalization**
 - Schema Matching – e.g., contact number and phone number
 - Compound attributes – full address vs str, city, state, zip
 - Nested attributes
 - List of features in one dataset (air conditioning, parking) vs feature a boolean attribute
 - Set valued attributes
 - Set of phones vs primary/secondary
 - Record segmentation from text
 - **Data normalization**
 - Often converting data to a standard form, remove whitespace
 - detecting and correcting errors in data values that contain known typographical errors or variations,
 - expanding abbreviations and replacing them with standard forms; replacing nicknames with their proper name forms
 - Usually done based on dictionaries (e.g., commercial dictionaries, postal addresses, etc.)
- Initial data prep big part of the work; smart normalization can go long way!

Matching Features

- **Comparison Vector γ :**

For two references x and y, compute a vector *similarity scores* of component attribute.

- [1st-author-match-score,
paper-match-score,
venue-match-score,
year-match-score,]

- Similarity scores

- Boolean (match or not-match)
- Real values based on distance functions
- Real values based on set or vector similarity

Summary of Matching Features

Permit efficient scalable
implementation

- Equality on a boolean predicate
 - Edit distance
 - Levenshtein, Smith-Waterman, Affine
 - Set similarity
 - Jaccard, Dice
 - Vector Based
 - Cosine similarity, TFIDF
 - Alignment-based or Two-tiered
 - Jaro-Winkler, Soft-TFIDF, Monge-Elkan
 - Phonetic Similarity
 - Soundex
 - Translation-based
 - Numeric distance between values
 - Domain-specific
-
- Useful packages:
 - SecondString, <http://secondstring.sourceforge.net/>
 - Simmetrics: <http://sourceforge.net/projects/simmetrics/>
 - LingPipe, <http://alias-i.com/lingpipe/index.html>

Summary of Matching Features

Handle Typographical errors

- Equality on a boolean predicate
- Edit distance
 - Levenshtein, Smith-Waterman, Affine
- Set similarity
 - Jaccard, Dice
- Vector Based
 - Cosine similarity, TFIDF

Good for Text (reviews/ tweets),
sets, class membership, ...

- Useful packages:
 - SecondString, <http://secondstring.sourceforge.net/>
 - Simmetrics: <http://sourceforge.net/projects/simmetrics/>
 - LingPipe, <http://alias-i.com/lingpipe/index.html>

Good for Names

- Alignment-based or Two-tiered
 - Jaro-Winkler, Soft-TFIDF, Monge-Elkan
- Phonetic Similarity
 - Soundex
- Translation-based
- Numeric distance between values
- Domain-specific

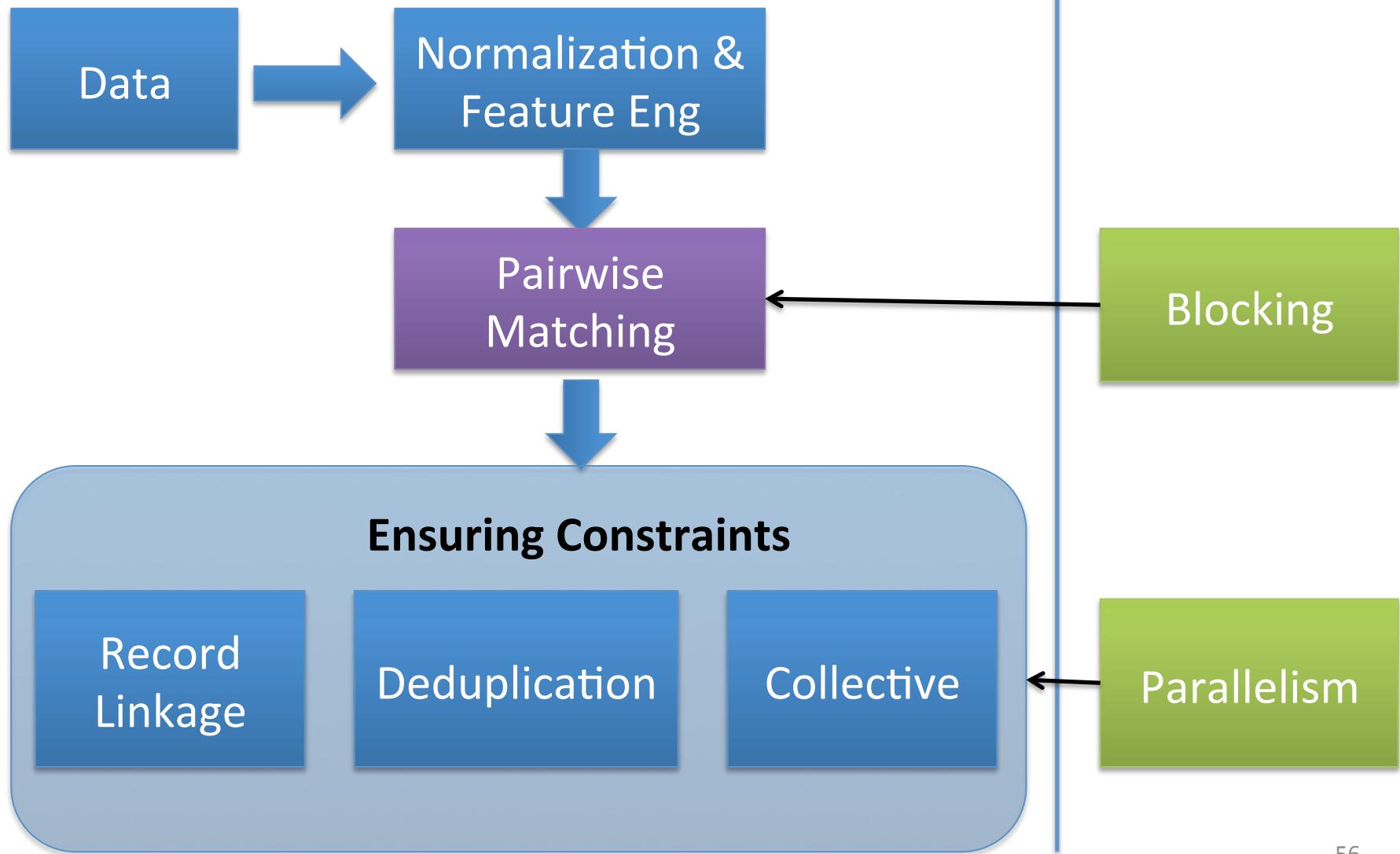
Useful for
abbreviations,
alternate names.

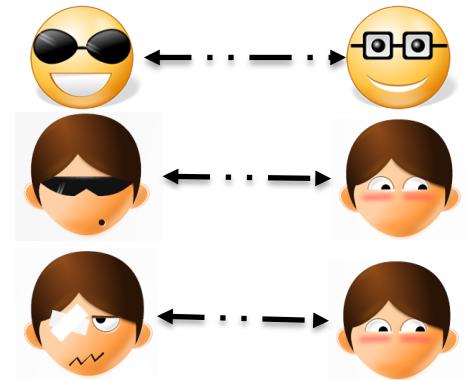
Relational Matching Features

- Relational features are often set-based
 - Set of coauthors for a paper
 - Set of cities in a country
 - Set of products manufactured by manufacturer
- Can use set similarity functions mentioned earlier
 - Common Neighbors: Intersection size
 - Jaccard's Coefficient: Normalize by union size
 - Adar Coefficient: Weighted set similarity
- Can reason about similarity in sets of values
 - Average or Max
 - Other aggregates

Part 2: Algorithmic Foundations

Part 3: Scaling





PART 2-b

PAIRWISE MATCHING

Pairwise Match Score

Problem: Given a vector of component-wise similarities for a pair of records (x, y) , compute $P(x \text{ and } y \text{ match})$.

Solutions:

1. Weighted sum or average of component-wise similarity scores.
Threshold determines match or non-match.
 - $0.5 * 1^{\text{st}}\text{-author-match-score} + 0.2 * \text{venue-match-score} + 0.3 * \text{paper-match-score}$.
 - Hard to pick weights.
 - Match on last name match *more predictive* than login name.
 - Match on “Smith” *less predictive* than match on “Getoor” or “Machanavajjhala”.
 - Hard to tune a threshold.

Pairwise Match Score

Problem: Given a vector of component-wise similarities for a pair of records (x, y) , compute $P(x \text{ and } y \text{ match})$.

Solutions:

1. Weighted sum or average of component-wise similarity scores.
Threshold determines match or non-match.
2. Formulate rules about what constitutes a match.
 - $(1^{\text{st}}\text{-author-match-score} > 0.7 \text{ AND venue-match-score} > 0.8)$
OR $(\text{paper-match-score} > 0.9 \text{ AND venue-match-score} > 0.9)$
 - Manually formulating the right set of rules is hard.

Fellegi & Sunter Model [FS, Science '69]

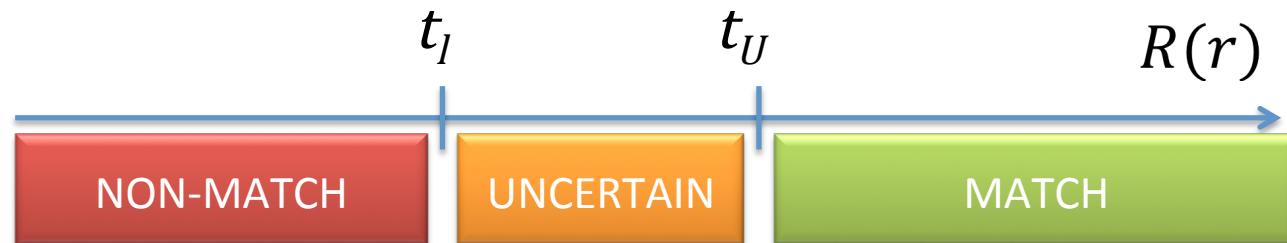
- Record pair: $r = (x, y)$ in $A \times B$
- $\gamma = \gamma(r)$ is a comparison vector
 - E.g., $\gamma = ["Is x.name = y.name?", "Is x.address = y.address?" ...]$
 - Assume binary vector for simplicity
- M : set of matching pairs of records
- U : set of non-matching pairs of records

Fellegi & Sunter Model [FS, Science '69]

- $r = (x, y)$ is record pair, γ is comparison vector, M matches, U non-matches
- Linkage decisions are based on:

$$R(r) = \frac{m(\gamma)}{u(\gamma)} = \frac{P(\gamma | r \in M)}{P(\gamma | r \in U)}$$

- **Linkage Rule: $L(t_l, t_u)$**



Error due to a Linkage Rule

- Type I Error: $r = (x,y)$ in U , but the linkage rule calls it a match

$$P(L_{match}|U) = \sum_{\gamma \in \Gamma} u(\gamma) \cdot P(L_{match}|\gamma)$$

- Type II Error: $r = (x,y)$ in M , but the linkage rule calls it a non-match

$$P(L_{non}|M) = \sum_{\gamma \in \Gamma} m(\gamma) \cdot P(L_{non}|\gamma)$$

Optimal Linkage Rule

- $L^* = (t_l^*, t_u^*)$ is an optimal decision rule for comparison space Γ with error bounds μ and λ , if

- L^* meets the type I and type II requirements

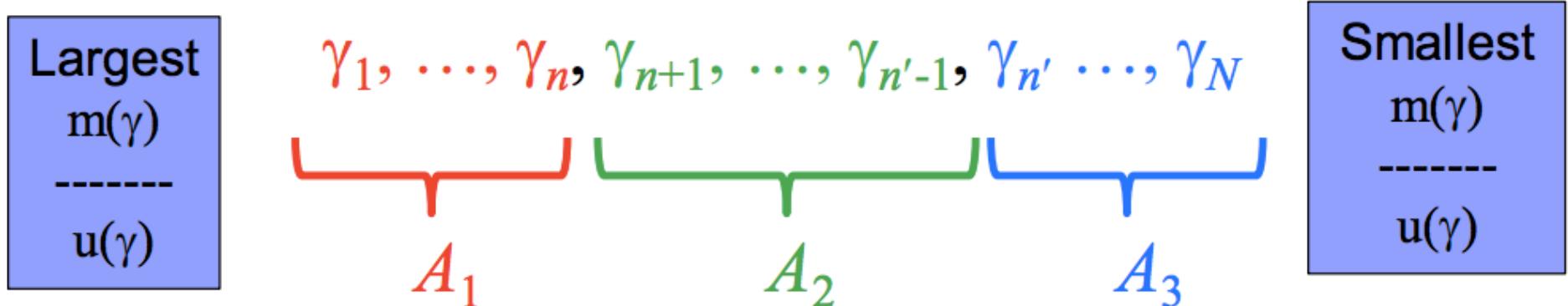
$$P(L_{match}|U) \leq \mu, \quad P(L_{non}|M) \leq \lambda$$

- L^* has the least conditional probabilities of *not making a decision*. That is for all other decision rules L (with error bounds μ and λ),

$$\begin{aligned} P(L_{uncertain}^*|U) &\leq P(L_{uncertain}|U) \\ P(L_{uncertain}^*|M) &\leq P(L_{uncertain}|M) \end{aligned}$$

Finding the Optimal Linkage Rule

- Suppose there are N comparison vectors
- Sort them in decreasing order of $m(\gamma) / u(\gamma)$



- Pick the largest n and n' such that:

$$\mu \geq \sum_{i=1}^n u(\gamma_i), \quad \lambda \geq \sum_{i=1}^n m(\gamma_i)$$

Using Fellegi Sunter in Practice

- Γ is usually high dimensional (computing $m(\gamma)$ and $u(\gamma)$ is inefficient)
 - Use conditional independence of features in γ given match or non-match
 - Naïve Bayes assumption
- Computing $P(\gamma \mid r \in M)$ requires some knowledge of matches.
 - Supervised learning (assume a training set is provided)
 - EM-based techniques can be used to learn the parameters jointly while identifying matches.

ML Pairwise Approaches

- Supervised machine learning algorithms
 - Decision trees
 - [Cochinwala et al, IS01]
 - Support vector machines
 - [Bilenko & Mooney, KDD03]; [Christen, KDD08]
 - Ensembles of classifiers
 - [Chen et al., SIGMOD09]
 - Conditional Random Fields (CRF)
 - [Gupta & Sarawagi, VLDB09]
 - ... and many others.
- Issues:
 - **Training set generation**
 - Imbalanced classes – many more negatives than positives (even after eliminating obvious non-matches ... using *Blocking*)
 - Misclassification cost

Creating a Training Set is a key issue

- Constructing a training set is hard – since most pairs of records are “easy non-matches”.
 - 100 records from 100 cities.
 - Only 10^6 pairs out of total 10^8 (1%) come from the same city
- Some pairs are hard to judge even by humans
 - Inherently ambiguous
 - E.g., Paris Hilton (person or business)
 - Missing attributes
 - Starbucks, Toronto vs Starbucks, Queen Street ,Toronto

Avoiding Training Set Generation

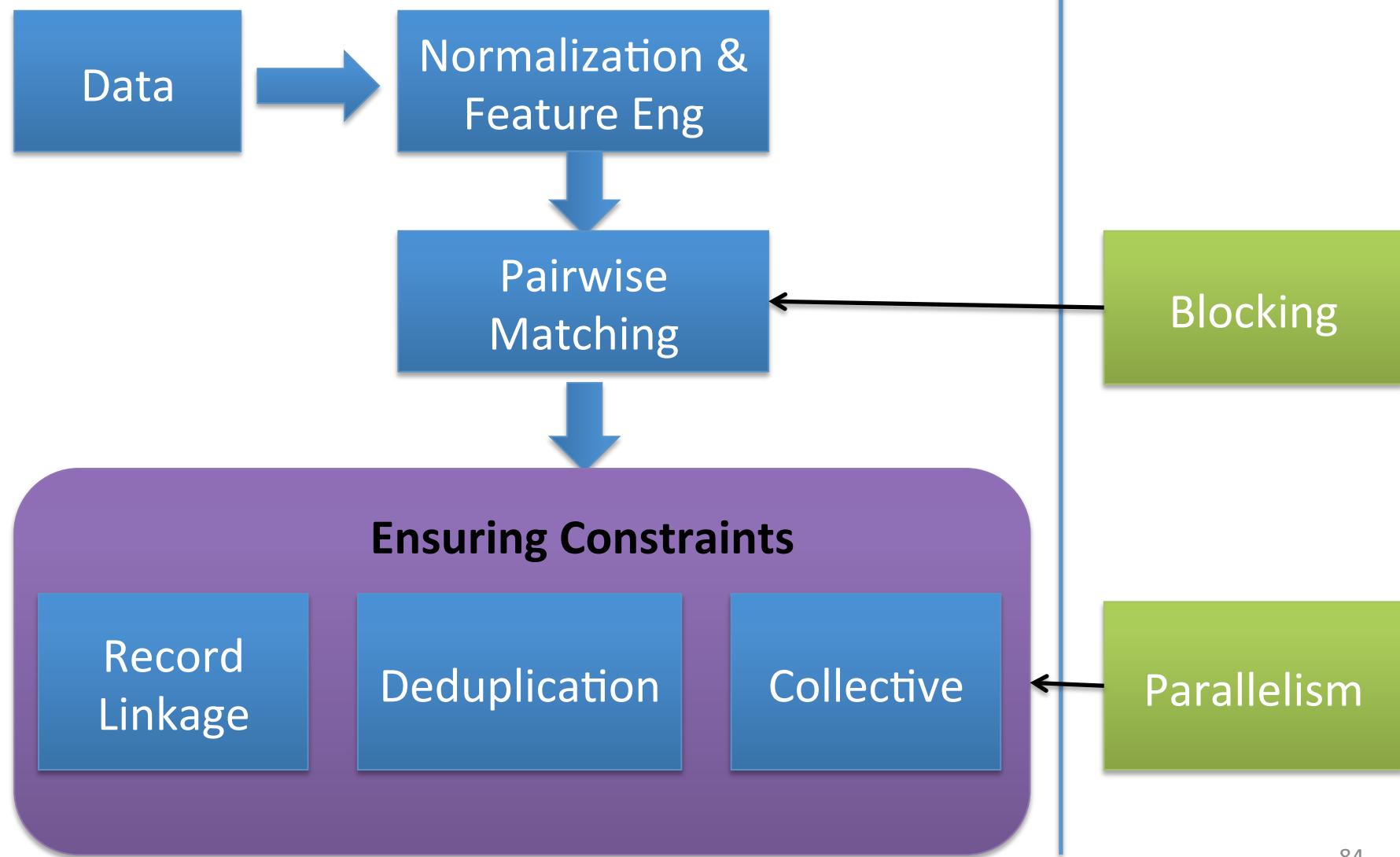
- Unsupervised / Semi-supervised Techniques
 - EM based techniques to learn parameters
 - [Winkler '06, Herzog et al '07]
 - Generative Models
 - [Ravikumar & Cohen, UAI04]
- Active Learning
 - Committee of Classifiers
 - [Sarawagi et al KDD '00, Tajeda et al IS '01]
 - Provably optimizing precision/recall
 - [Arasu et al SIGMOD '10, Bellare et al KDD '12]
 - Crowdsourcing
 - [Wang et al VLDB '12, Marcus et al VLDB '12, ...]

Summary of Pairwise Matching

- Many algorithms for independent classification of pairs of records as match/non-match
- ML based classification & Fellegi-Sunter
 - Pro: Advanced state of the art
 - Con: Building high fidelity training sets is a hard problem
- Active Learning & Crowdsourcing for ER are active areas of research.

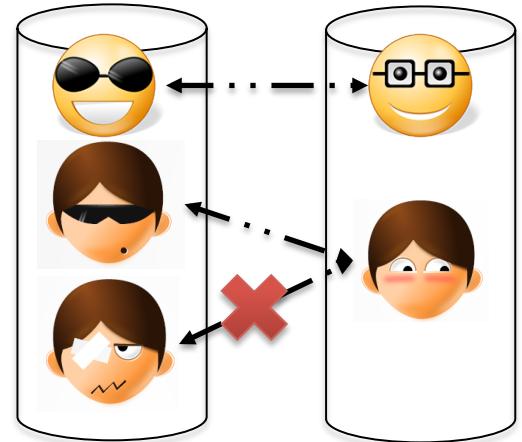
Part 2: Algorithmic Foundations

Part 3: Scaling



Algorithms for Handling Constraints

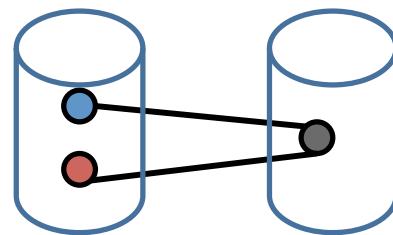
- Record linkage - propagation through exclusivity
 - Weighted k-partite matching
- Deduplication - propagation through transitivity
 - Correlation clustering
- Collective - propagation through general constraints
 - Similarity propagation
 - Dependency graphs, Collective Relational Clustering
 - Probabilistic approaches
 - LDA, CRFs, Markov Logic Networks, Probabilistic Relational Models,
 - Hybrid approaches
 - Dedupalog



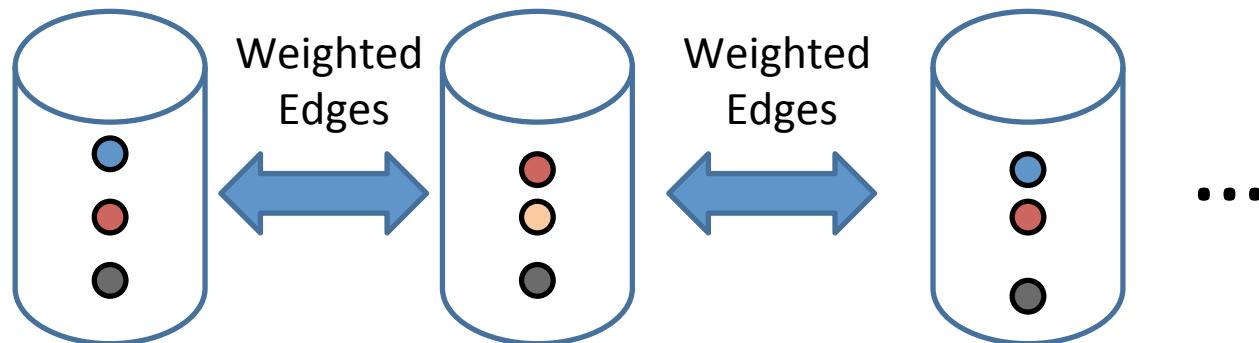
RECORD LINKAGE

1-1 assumption

- Matching between (almost) deduplicated databases.
- Each record in one database matches at most one record in another database.
- Pairwise ER may match a record in one database with more than one record in second database

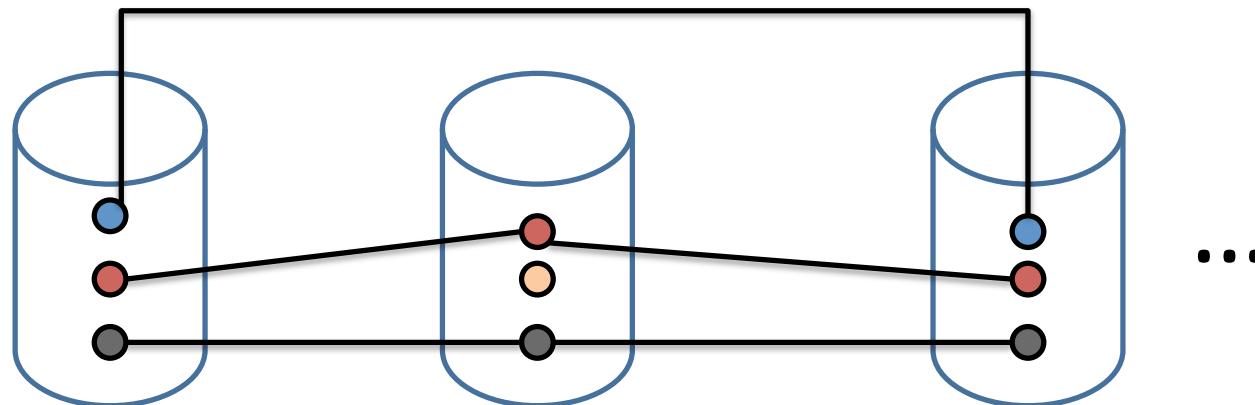


Weighted K-Partite Matching



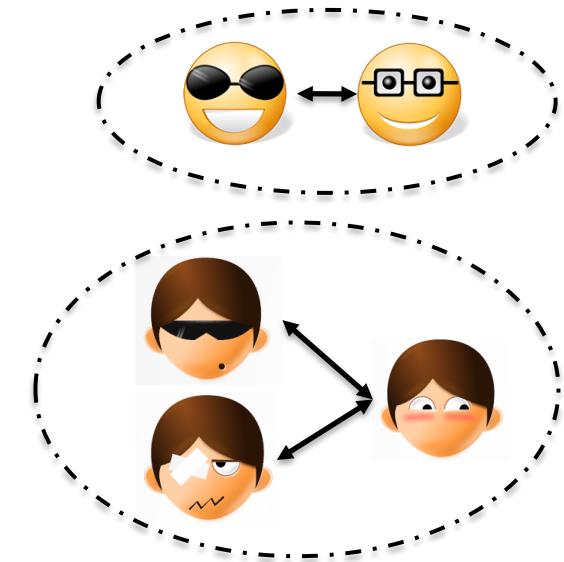
- Edges between pairs of records from different databases
- Edge weights
 - Pairwise match score
 - Log odds of matching

Weighted K-Partite Matching



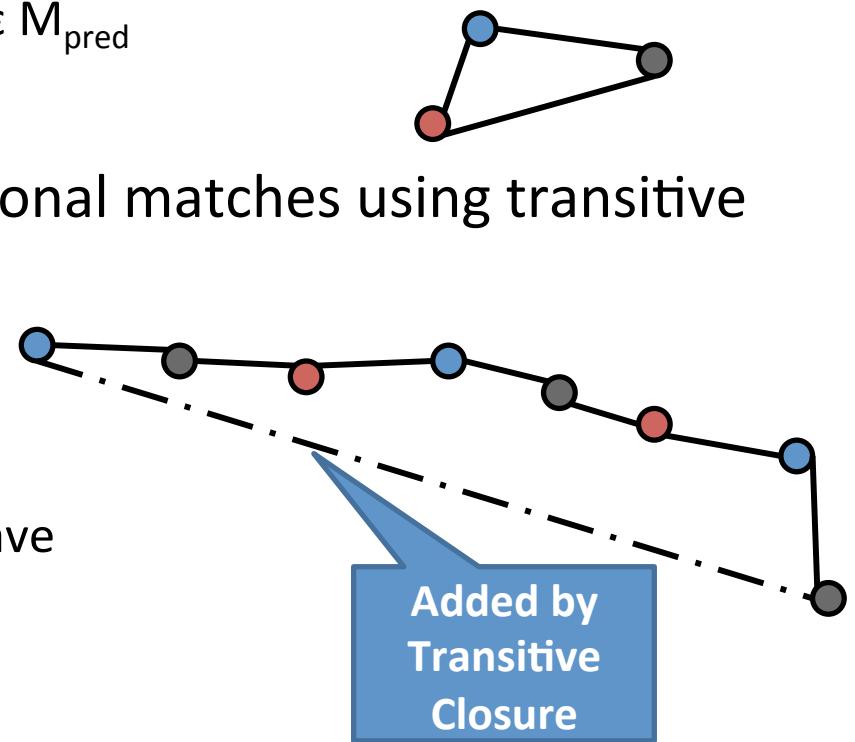
- Find a matching (each record matches at most one other record from other database) that maximize the sum of weights.
- General problem is NP-hard (3D matching)
- Successive bipartite matching is typically used. [Gupta & Sarawagi, VLDB '09]

DEDUPLICATION



Deduplication => Transitivity

- Often pairwise ER algorithm output “inconsistent” results
 - $(x, y) \in M_{\text{pred}}$, $(y, z) \in M_{\text{pred}}$, but $(x, z) \notin M_{\text{pred}}$
- Idea: Correct this by adding additional matches using transitive closure
- In certain cases, this is a bad idea.
 - Graphs resulting from pairwise ER have diameter > 20
[Rastogi et al ICDE ‘13]
- Need clustering solutions that deal with this problem directly by reasoning about records jointly.



Clustering-based ER

- Resolution decisions are not made independently for each pair of records
- Based on variety of clustering algorithms, but
 - Number of clusters unknown aprioiri
 - Many, many small (possibly singleton) clusters
- Often take a pair-wise similarity graph as input
- May require the construction of a *cluster representative* or *canonical entity*

Clustering Methods for ER

- Hierarchical Clustering
 - [Bilenko et al, ICDM 05]
- Nearest Neighbor based methods
 - [Chaudhuri et al, ICDE 05]
- **Correlation Clustering**
 - [Soon et al CL'01, Bansal et al ML'04, Ng et al ACL'02, Ailon et al JACM'08, Elsner et al ACL'08, Elsner et al ILP-NLP'09]

Integer Linear Programming view of ER

- $r_{xy} \in \{0,1\}$, $r_{xy} = 1$ if records x and y are in the same cluster.
- $w^+_{xy} \in [0,1]$, benefit of clustering x and y together
- $w^-_{xy} \in [0,1]$, benefit of placing x and y in different clusters

$$\text{maximize } \sum r_{xy} w^+_{xy} + (1 - r_{xy}) w^-_{xy}$$

s.t. $\forall x, y, z \in R,$

$$r_{xy} + r_{xz} + r_{yz} \neq 2$$

Transitive
closure

Correlation Clustering

$$\text{maximize} \sum r_{xy} w_{xy}^+ + (1 - r_{xy}) w_{xy}^-$$

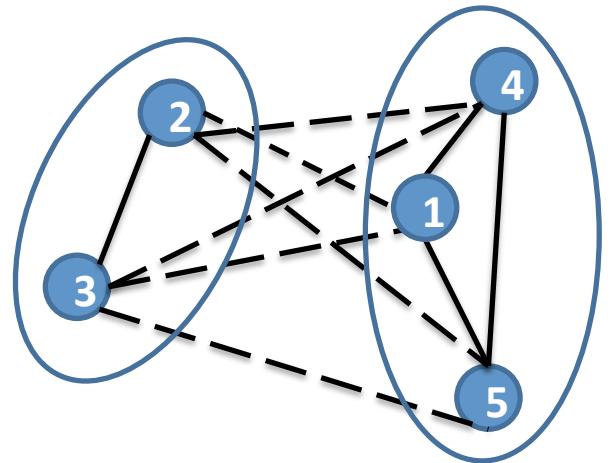
s.t. $\forall x, y, z \in R,$

$$r_{xy} + r_{xz} + r_{yz} \neq 2$$

- Cluster mentions such that total benefit is maximized

Solid edges contribute w_{xy}^+ to the objective

Dashed edges contribute w_{xy}^- to the objective



- Benefit based on pairwise similarities

$$\{p_{xy} \mid \forall (x, y) \in R \times R\}$$

- Additive: $w_{xy}^+ = p_{xy}$ and $w_{xy}^- = (1-p_{xy})$
- Logarithmic: $w_{xy}^+ = \log(p_{xy})$ and $w_{xy}^- = \log(1-p_{xy})$

Correlation Clustering

- Solving the ILP is NP-hard [Ailon et al 2008 JACM]
- A number of heuristics [Elsner et al 2009 ILP-NLP]
 - Greedy BEST/FIRST/VOTE algorithms
 - Greedy PIVOT algorithm (5-approximation)
 - Local Search

Greedy Algorithms

Step 1: Permute the nodes according a random π

Step 2: Assign record x to the cluster that maximizes *Quality*
Start a new cluster if $Quality < 0$

Quality:

- BEST: Cluster containing the closest match $\max_{y \in C} w_{xy}^+$
 - [Ng et al 2002 ACL]
- FIRST: Cluster contains the most recent vertex y with $w_{xy}^+ > 0$
 - [Soon et al 2001 CL]
- VOTE: Assign to cluster that minimizes objective function.
 - [Elsner et al 08 ACL]

Practical Note:

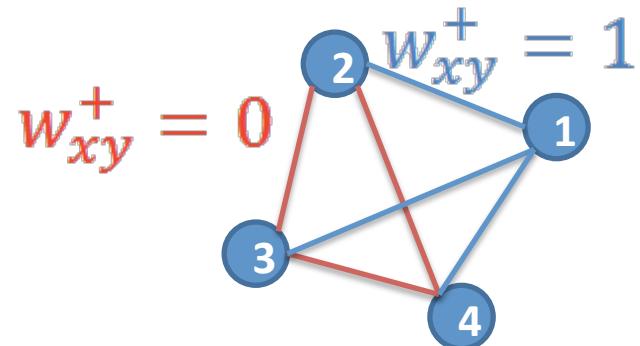
- Run the algorithm for many random permutations , and pick the clustering with best objective value (better than average run)

Greedy with approximation guarantees

PIVOT Algorithm

[Ailon et al 2008 JACM]

- Pick a random (*pivot*) record p .
- New cluster = $\{x \mid w_{px}^+ > 0\}$
- $\pi = \{1,2,3,4\}$ $C = \{\{1,2,3,4\}\}$
- $\pi = \{2,4,1,3\}$ $C = \{\{1,2\}, \{4\}, \{3\}\}$
- $\pi = \{3,2,4,1\}$ $C = \{\{1,3\}, \{2\}, \{4\}\}$



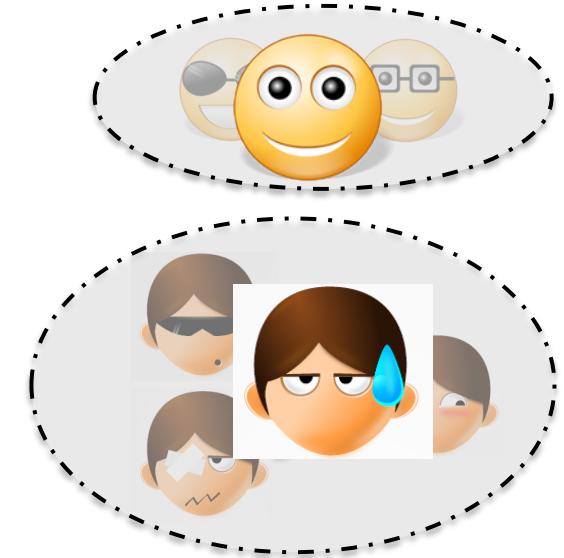
When weights are 0/1,

$E(\text{cost(greedy)}) < 3 \text{ OPT}$

For $w_{xy}^+ + w_{xy}^- = 1$,

$E(\text{cost(greedy)}) < 5 \text{ OPT}$

[Elsner et al, ILP-NLP '09] : Comparison of various correlation clustering algorithms



CANONICALIZATION

Canonicalization

- Merge information from duplicate mentions to construct a cluster representative with *maximal* information

- Starbucks,
3457 Hillsborough Road
Durham, NC
Ph: *null*
- Starbacks,
Hillsborough Rd, Durham
Ph: (919) 333-4444

Starbucks
3457 Hillsborough Road, Durham, NC
Ph: (919) 333-4444

Critically important in Web portals where users must be shown a consolidated view

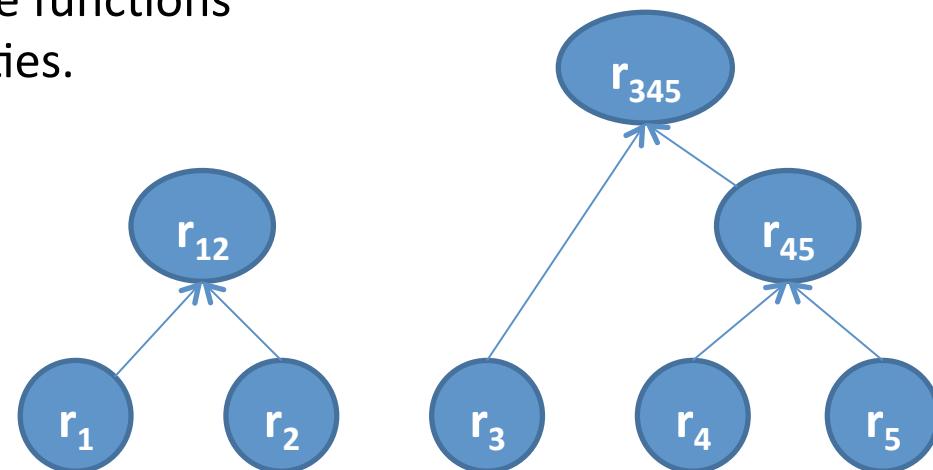
- Each mention only contains a subset of the attributes
- Mentions contain variations (of names, addresses)
- Some of the mentions have incorrect values

Canonicalization Algorithms

- Rule based:
 - For names: typically longest names are used.
 - For set values attributes: UNION is used.
- For strings, [Culotta et al KDD07] learn an edit distance for finding the most representative “centroid”.
- Can use “majority rule” to fix errors
(if 4 out of 5 say a business is closed, then business is closed).
 - This may not always work due to copying [Dong et al VLDB09], or when underlying data changes [Pal et al WWW11]

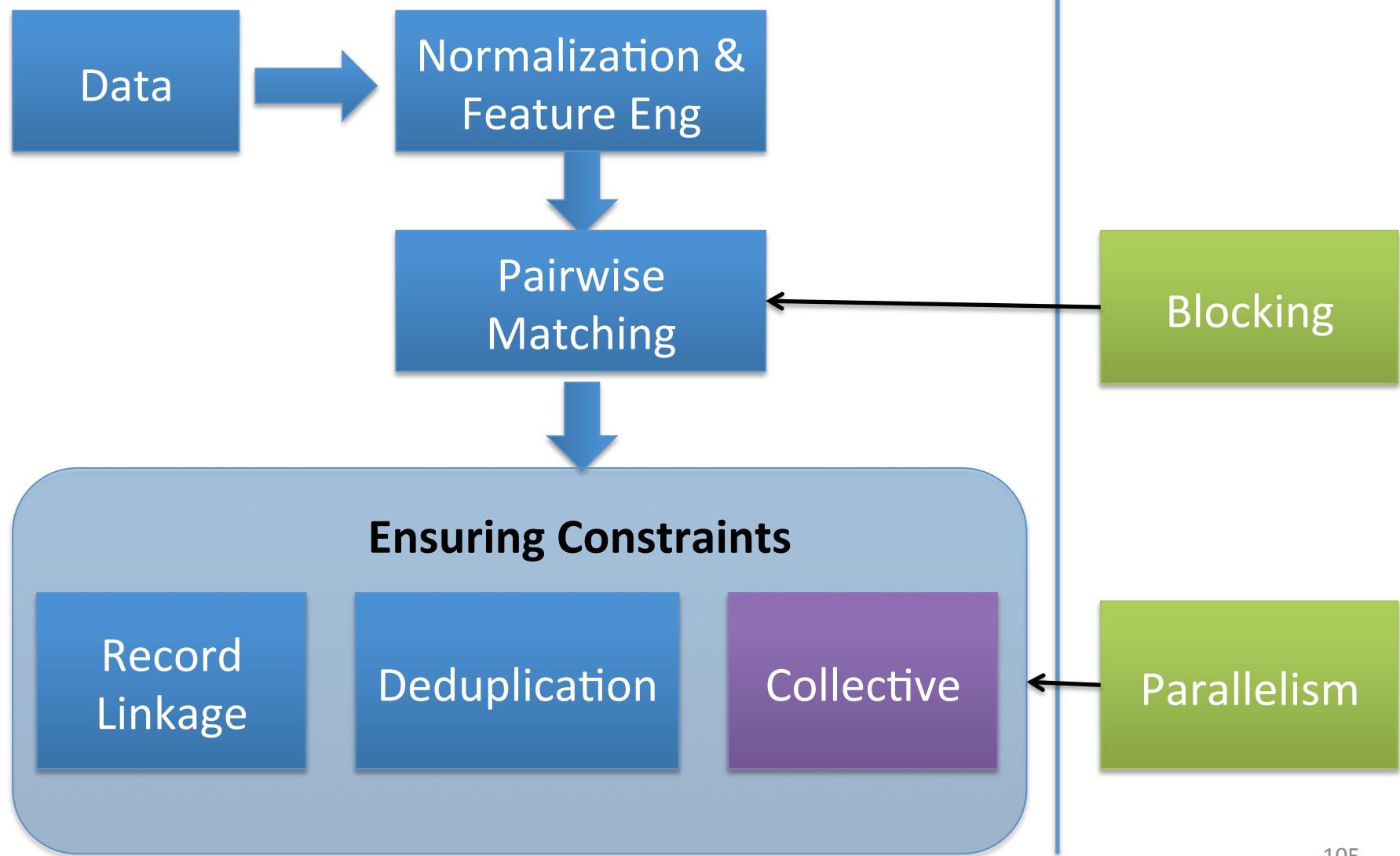
Canonicalization for Efficiency

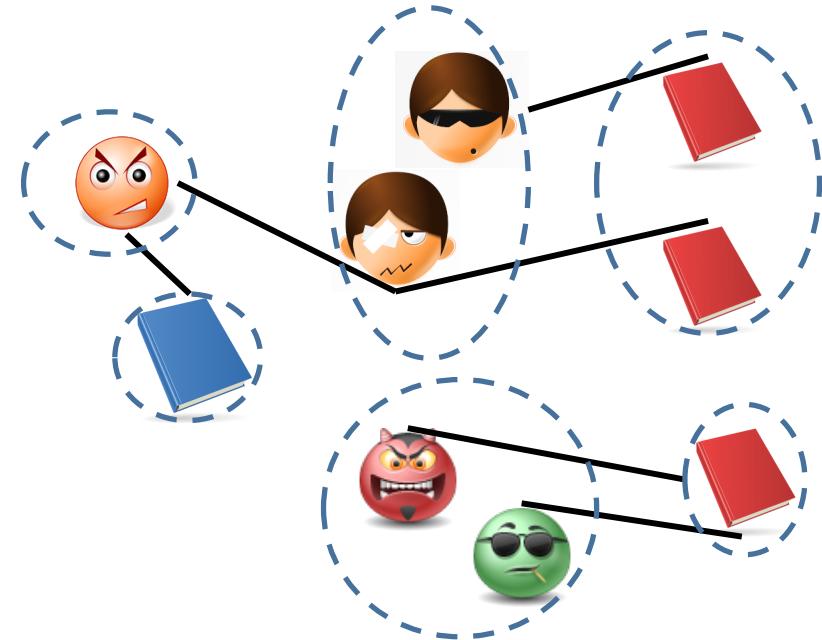
- Stanford Entity Resolution Framework [Benjelloun VLDBJ09]
 - Consider a blackbox match and merge function
 - Match is a pairwise boolean operator
 - Merge: construct canonical version of a matching pair
- Can minimize time to compute matches by interleaving matching and merging
 - esp., when match and merge functions satisfy **monotonicity** properties.



Part 2: Algorithmic Foundations

Part 3: Scaling

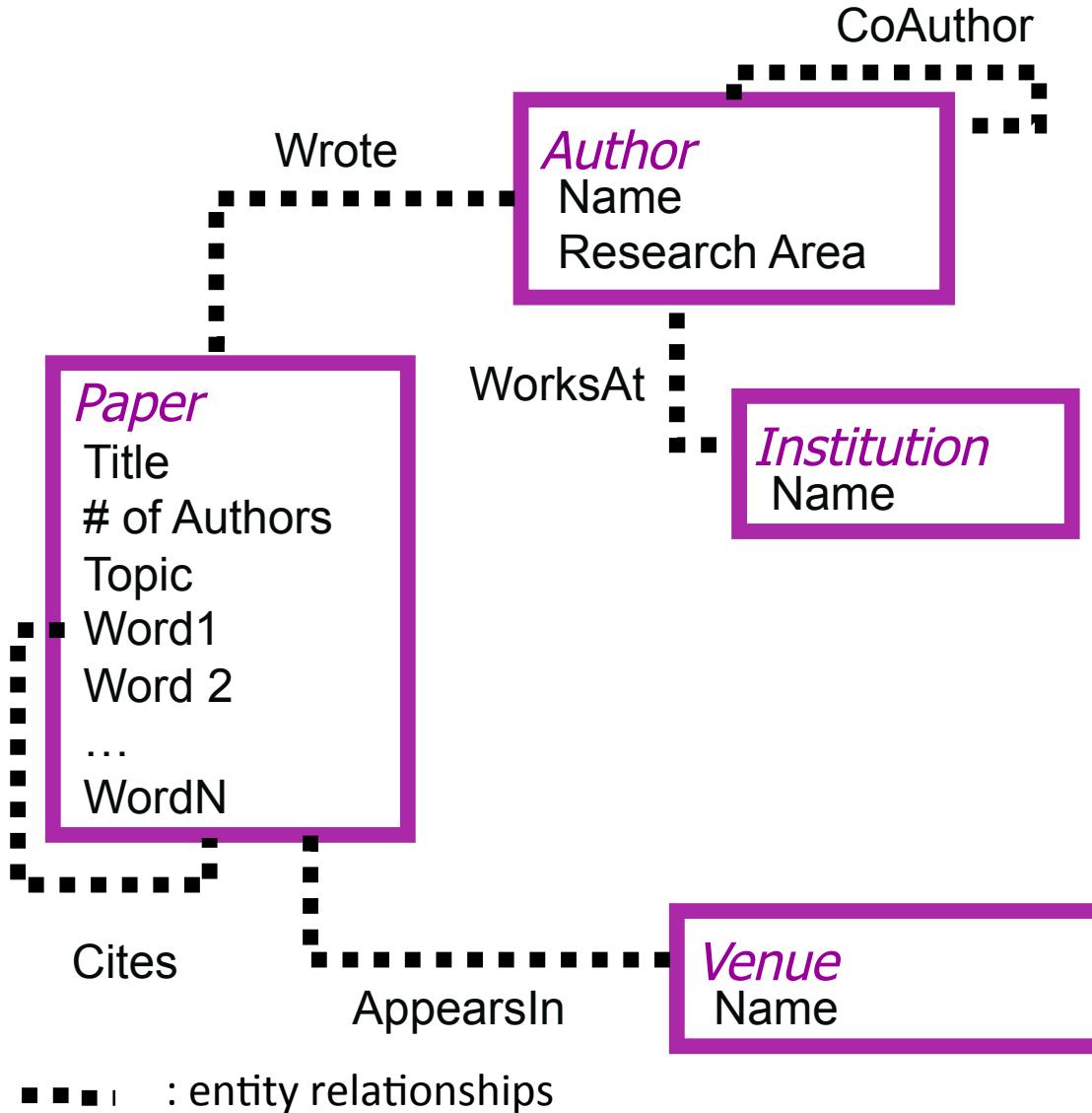




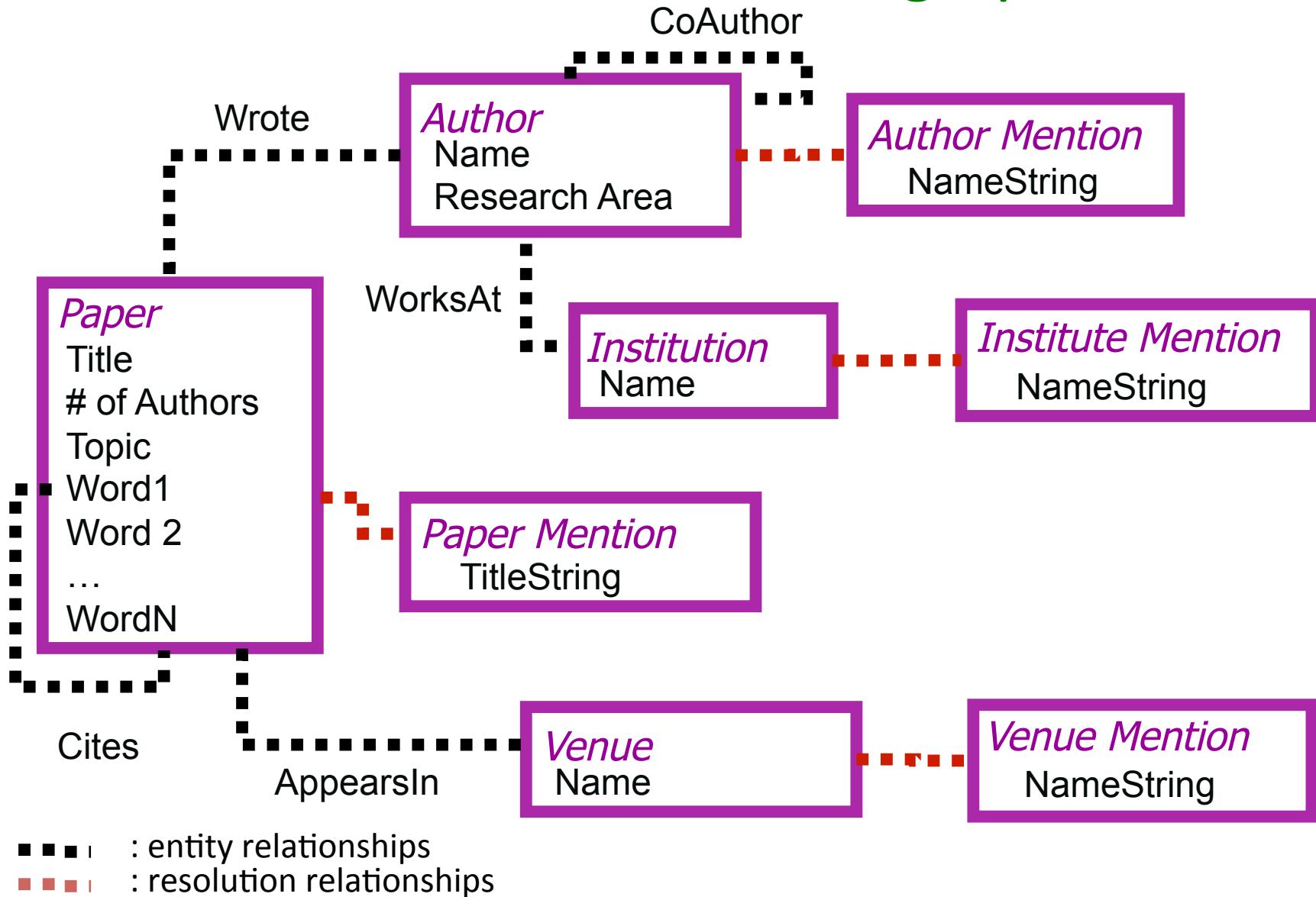
COLLECTIVE ENTITY RESOLUTION

MOTIVATING EXAMPLE: BIBLIOGRAPHIC DOMAIN

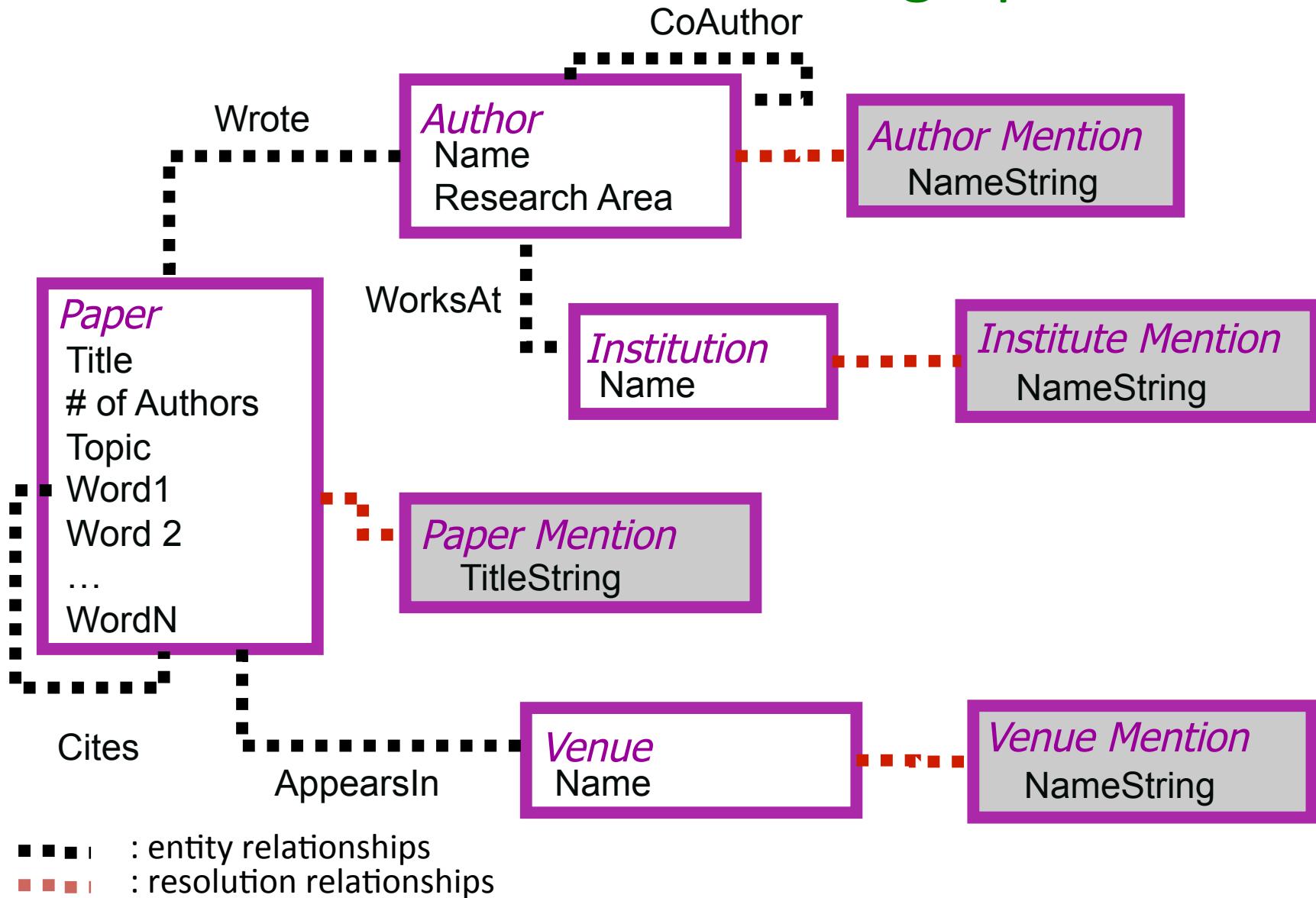
Entities & Relations in Bibliographic Domain



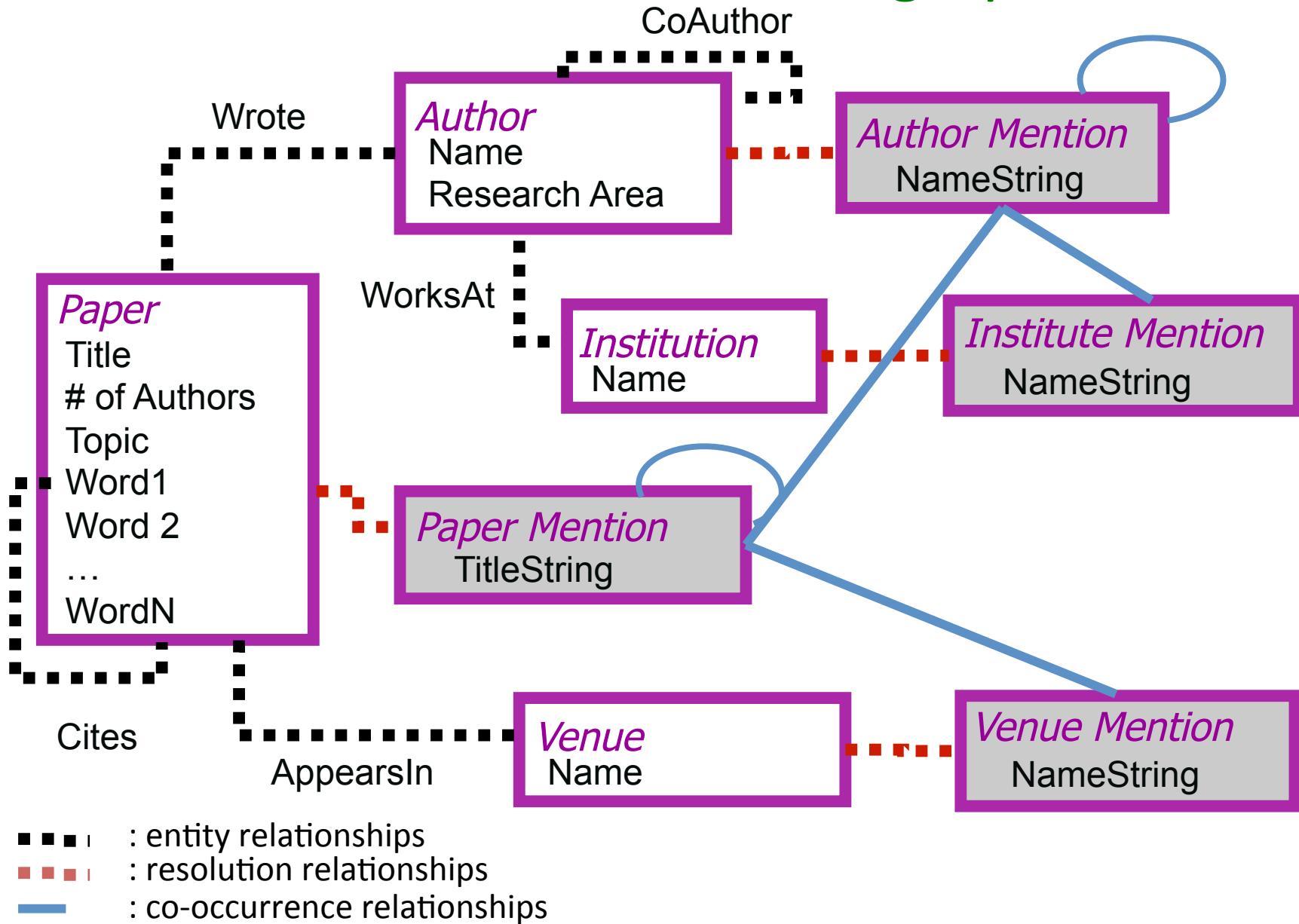
Entities & Relations in Bibliographic Domain



Entities & Relations in Bibliographic Domain



Entities & Relations in Bibliographic Domain



Constraints

- Basic constraints (review):
 - **Transitivity:** If M1 and M2 match, M2 and M3 match, then M1 and M3 match
 - **Exclusivity:** If M1 matches with M2, then M3 cannot match with M2
 - **Functional Dependency:** If M1 and M2 match, then M3 and M4 must match
- General constraints:
 - Positive/Negative evidence for match
 - Hard/soft constraints
- Collective Entity Resolution methods
 - Resolutions decisions are not independent!

Constraint Types

	Hard Constraint	Soft Constraint
Positive Evidence		
Negative Evidence		

Constraint Types

	Hard Constraint	Soft Constraint
Positive Evidence	<p>If M1, M2 match then M3, M4 must match (functional dependency)</p> <p><i>If two papers match, their venues match</i></p> <p><i>If two papers match, their set of coauthors match</i></p>	
Negative Evidence		

Constraint Types

	Hard Constraint	Soft Constraint
Positive Evidence	<p>If M1, M2 match then M3, M4 must match (functional dependency)</p> <p><i>If two papers match, their venues match</i></p> <p><i>If two papers match, their set of coauthors match</i></p>	
Negative Evidence	<p>Mention M1 and M2 must refer to distinct entities (Uniqueness)</p> <p><i>Coauthors are distinct</i></p> <p>If M1, M2 don't match then M3, M4 cannot match</p> <p><i>If two venues don't match, then their papers don't match</i></p>	

Constraint Types

	Hard Constraint	Soft Constraint
Positive Evidence	<p>If M1, M2 match then M3, M4 must match (functional dependency)</p> <p><i>If two papers match, their venues match</i></p> <p><i>If two papers match, their set of coauthors match</i></p>	<p>If M1, M2 match then M3, M4 more likely to match</p> <p><i>If two venues match, then their papers are more likely to match</i></p> <p><i>If two authors are a match, their coauthors are more likely to match</i></p> <p><i>If two authors are from the same research group, they are more likely to be coauthors</i></p>
Negative Evidence	<p>Mention M1 and M2 must refer to distinct entities (Uniqueness)</p> <p><i>Coauthors are distinct</i></p> <p>If M1, M2 don't match then M3, M4 cannot match</p> <p><i>If two venues don't match, then their papers don't match</i></p>	

Constraint Types

	Hard Constraint	Soft Constraint
Positive Evidence	<p>If M1, M2 match then M3, M4 must match (functional dependency)</p> <p><i>If two papers match, their venues match</i></p> <p><i>If two papers match, their set of coauthors match</i></p>	<p>If M1, M2 match then M3, M4 more likely to match</p> <p><i>If two venues match, then their papers are more likely to match</i></p> <p><i>If two authors are a match, their coauthors are more likely to match</i></p> <p><i>If two authors are from the same research group, they are more likely to be coauthors</i></p>
Negative Evidence	<p>Mention M1 and M2 must refer to distinct entities (Uniqueness)</p> <p><i>Coauthors are distinct</i></p> <p>If M1, M2 don't match then M3, M4 cannot match</p> <p><i>If two venues don't match, then their papers don't match</i></p>	<p>If M1, M2 don't match then M3, M4 less likely to match</p> <p><i>If institutions don't match, then authors less likely to match</i></p>

Constraint Types

	Hard Constraint	Soft Constraint
Positive Evidence	<p>If M1, M2 match then M3, M4 must match (functional dependency)</p> <p><i>If two papers match, their venues match</i></p> <p><i>If two papers match, their set of coauthors may be</i> set-based</p> <p>constraints may be relational and require joins</p>	<p>If M1, M2 match then M3, M4 more likely to match</p> <p><i>If two venues match, then their papers are more likely to match</i></p> <p><i>If two authors are a match, their coauthors are more likely to match</i></p> <p><i>If two authors are from the same research group, they are more likely to be coauthors</i></p> <p>recursive</p> <p>latent variable</p>
Negative Evidence	<p>Mention M1 and M2 must refer to distinct entities (Uniqueness)</p> <p><i>Coauthors are distinct</i></p> <p>If M1, M2 don't match then M3, M4 cannot match</p> <p><i>If two venues don't match, then their papers don't match</i></p>	<p>If M1, M2 don't match then M3, M4 less likely to match</p> <p><i>If institutions don't match, then authors less likely to match</i></p>

Match Extent

- **Global:** If two papers match, then their venues match
 - This constraint can be applied to all instances of venue mentions
 - All occurrences of ‘SIGMOD’ can be matched to ‘International Conference on Management of Data’
- **Local:** If two papers match, then their authors match
 - This constraint can only be applied locally
 - Don’t want to match all occurrences of ‘J. Smith’ with ‘Jeff Smith’, only in the context of the current paper

Additional Constraints

Type	Example
Aggregate	C1 = No researcher has published more than five AAAI papers in a year
Subsumption	C2 = If a citation X from DBLP matches a citation Y in a homepage, then each author mentioned in Y matches some author mentioned in X
Neighborhood	C3 = If authors X and Y share similar names and some co-authors, they are likely to match
Incompatible	C4 = No researcher exists who has published in both HCI and numerical analysis
Layout	C5 = If two mentions in the same document share similar names, they are likely to match
Key/Uniqueness	C6 = Mentions in the PC listing of a conference is to different researchers
Ordering	C7 = If two citations match, then their authors will be matched in order
Individual	C8 = The researcher with the name “Mayssam Saria” has fewer than five mentions in DBLP (new graduate student)

Match Dependencies

When matching decisions depend on other matching decisions (in other words, matching decisions are not made independently), we refer to the approach as ***collective***

Collective Approaches

- Decisions for cluster-membership depends on other clusters
 - Non-probabilistic approaches
 - Similarity Propagation
 - Multi-relational Clustering Approaches
 - Probabilistic Models
 - Generative Models
 - Undirected Models
 - Hybrid Approaches

SIMILARITY PROPAGATION

Similarity Propagation Approaches

- Similarity propagation algorithms define a graph which encodes the similarity between entity mentions and matching decisions, and compute matching decisions by propagating similarity values.
 - Details of constructed graph and how the similarity is computed varies
 - Algorithms are usually defined procedurally
 - While probabilities may be encoded in various ways in the algorithms, no global probabilistic model is defined
- Approaches often more scalable than global probabilistic models

Dependency Graph

[Dong et al., SIGMOD05]

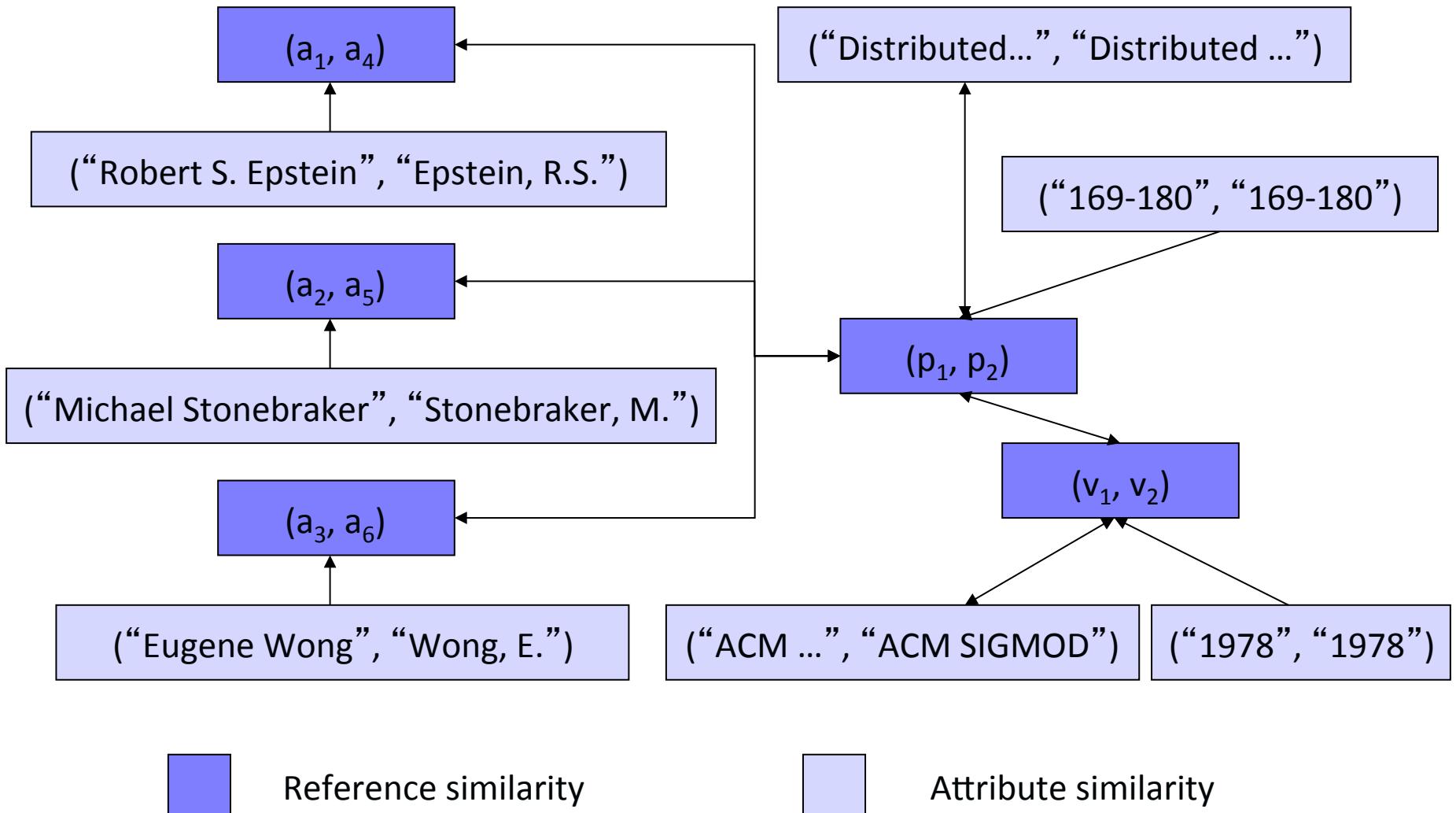
- Construct a graph where nodes represent similarity comparisons between attribute values (real-valued) and match decisions based on matching decisions of associated nodes (boolean-valued)
- As mentions are resolved, enriched to contain associated nodes of all matched mentions
- Similarity propagated until fixed point is reached
- Negative constraints (not-match nodes) are checked after similarity propagation is performed, and inconsistencies are fixed

Real-World Objects

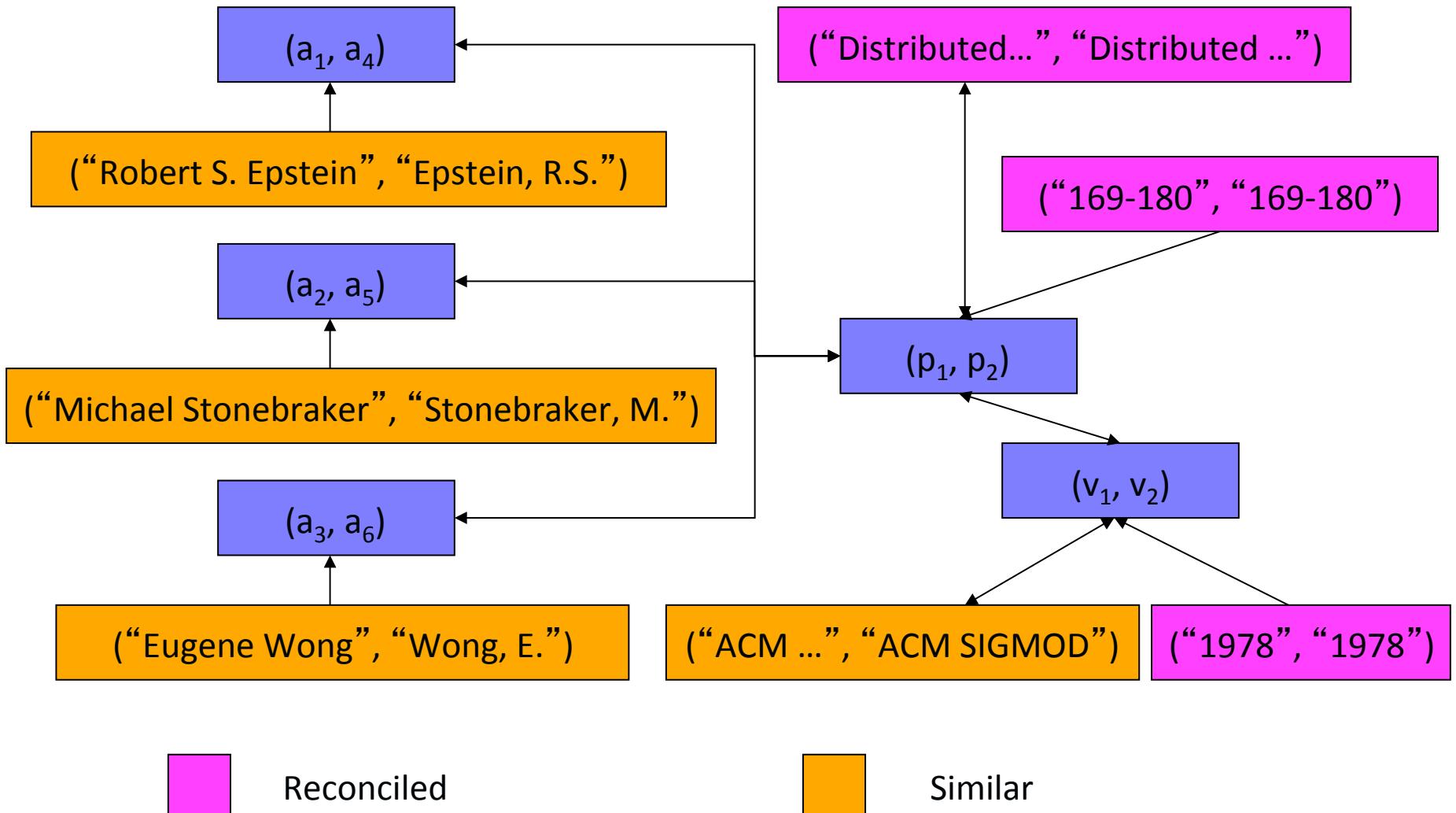
- **Paper:** $p_1 = ("Distributed\ Query\ Processing", "169-180", \{a_1, a_2, a_3\}, v_1)$
 $p_2 = ("Distributed\ query\ processing", "169-180", \{a_4, a_5, a_6\}, v_2)$
- **Venue:** $v_1 = ("ACM\ Conference\ on\ Management\ of\ Data", "1978",$
“Austin, Texas”)
 $v_2 = ("ACM\ SIGMOD", "1978", null)$
- **Author:** $a_1 = ("Robert\ S.\ Epstein", null)$
 $a_2 = ("Michael\ Stonebraker", null)$
 $a_3 = ("Eugene\ Wong", null)$
 $a_4 = ("Epstein,\ R.S.", null)$
 $a_5 = ("Stonebraker,\ M.", null)$
 $a_6 = ("Wong,\ E.", null)$

[Based on slides from Luna Dong]

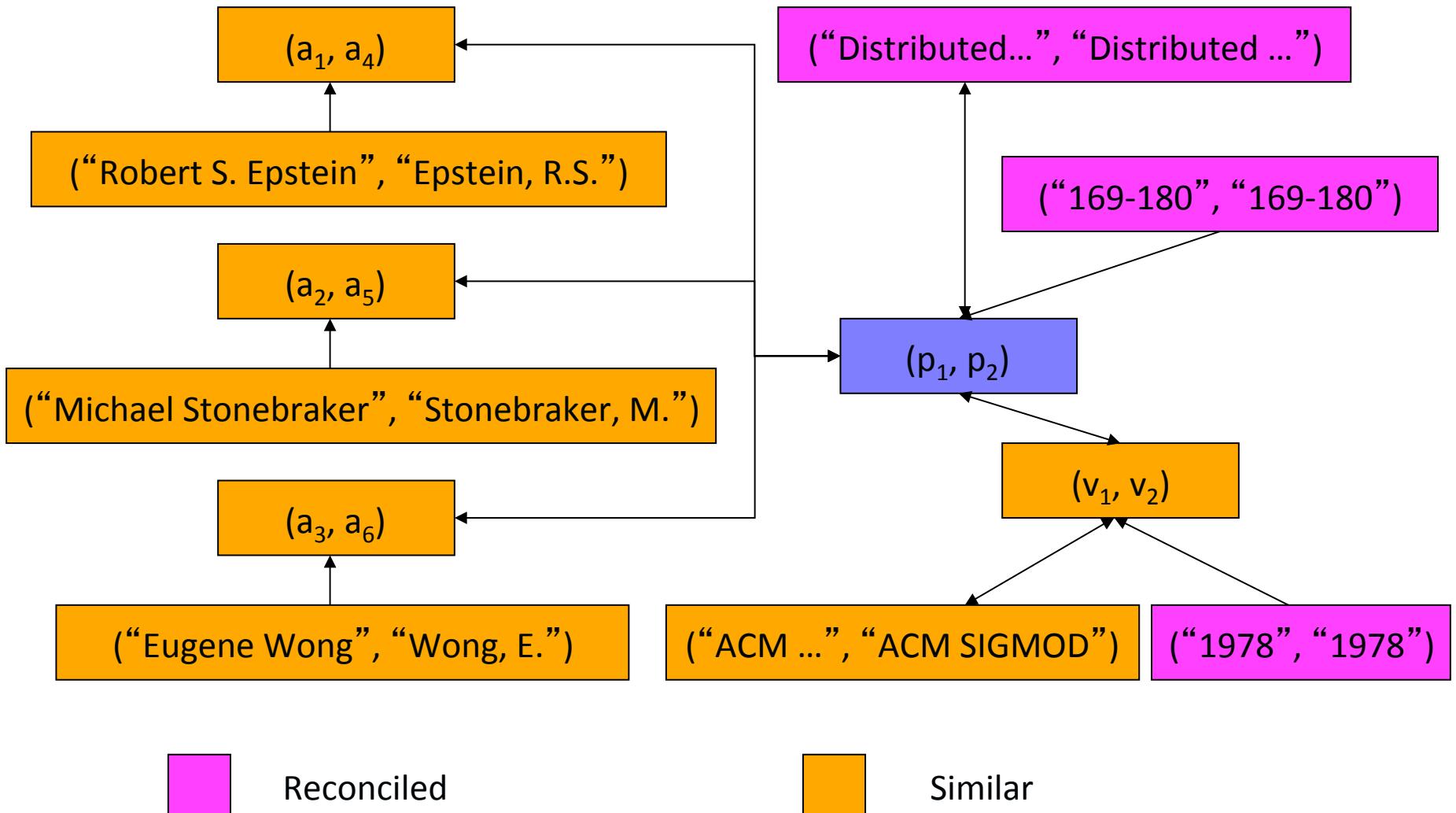
Exploit the Dependency Graph



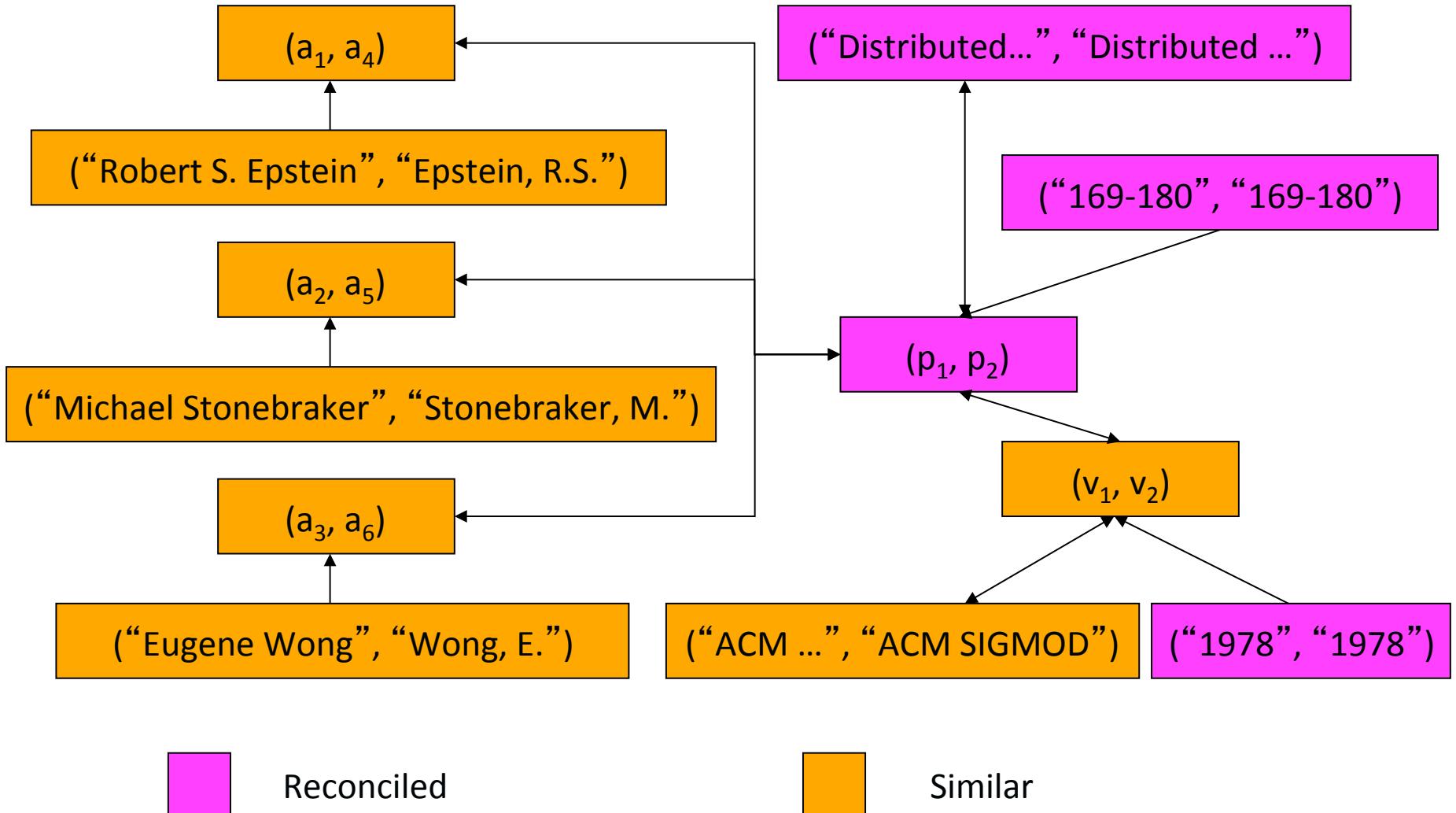
Exploit the Dependency Graph



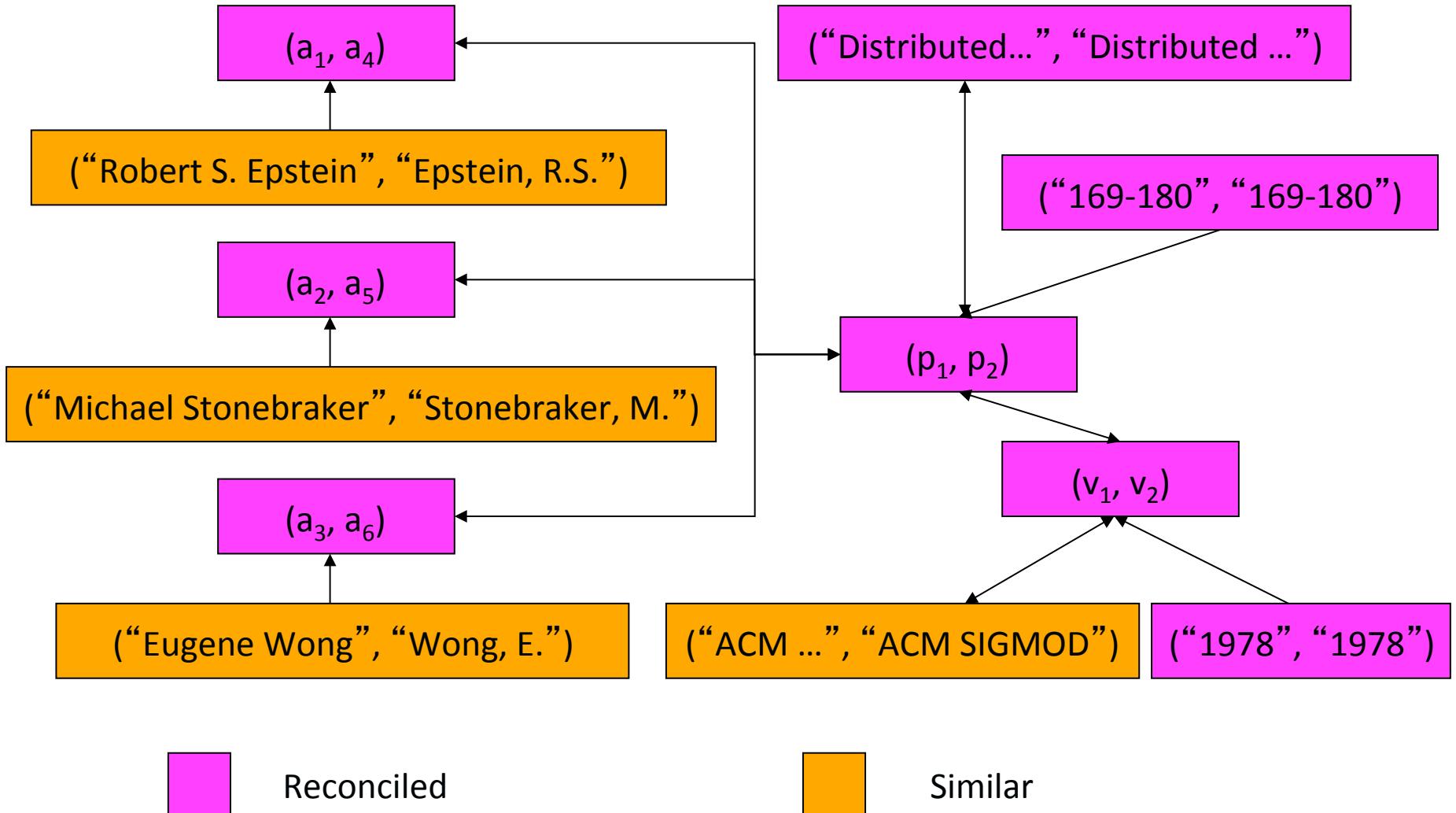
Exploit the Dependency Graph



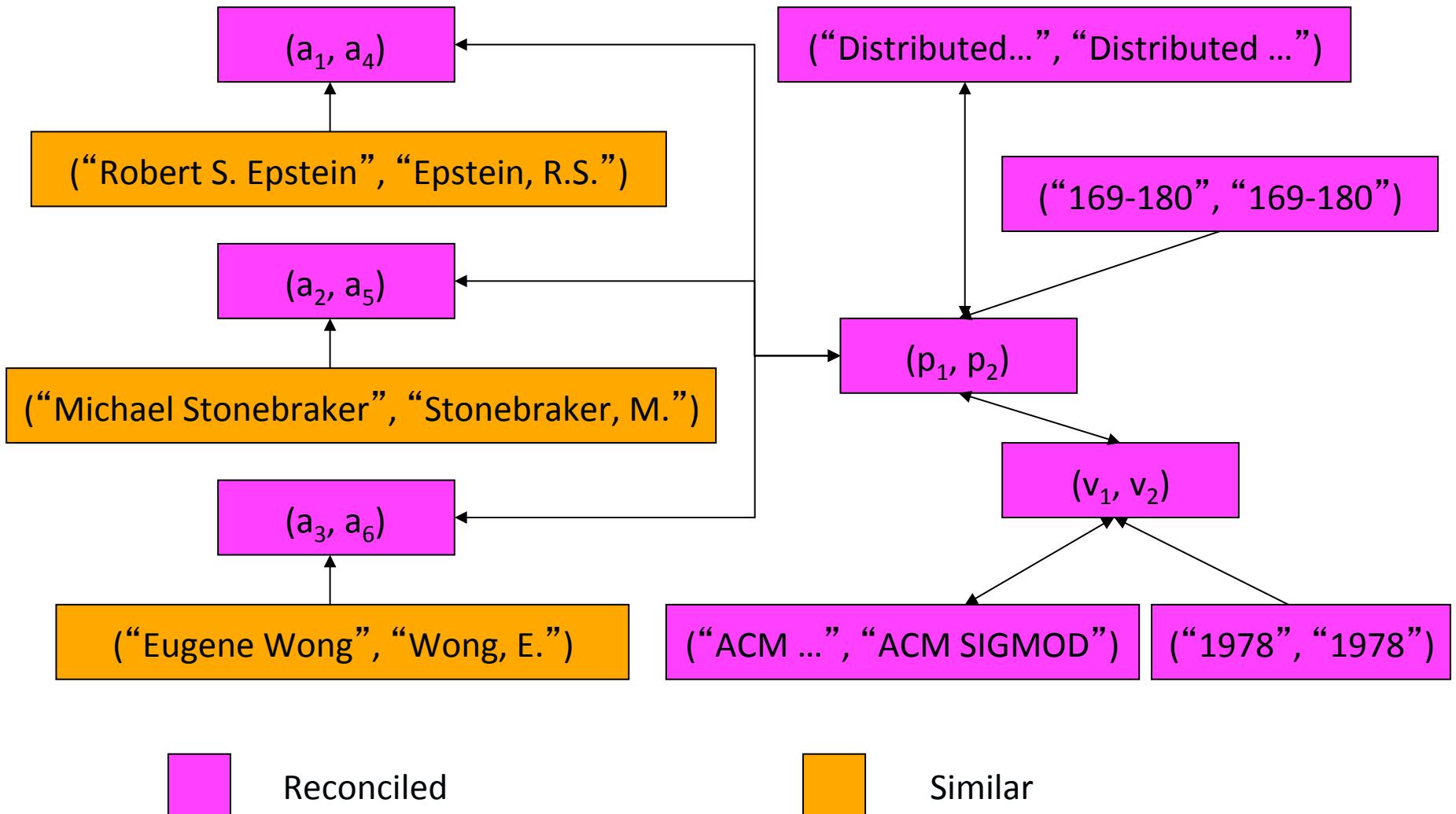
Exploit the Dependency Graph



Exploit the Dependency Graph

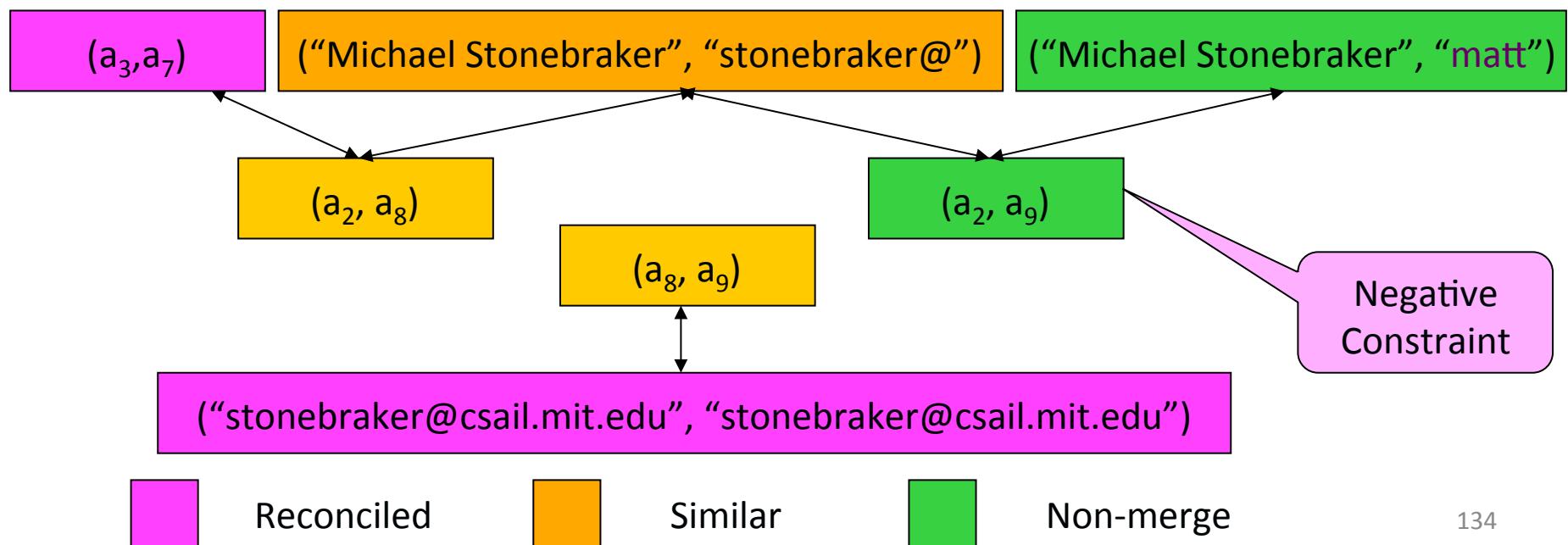


Exploit the Dependency Graph



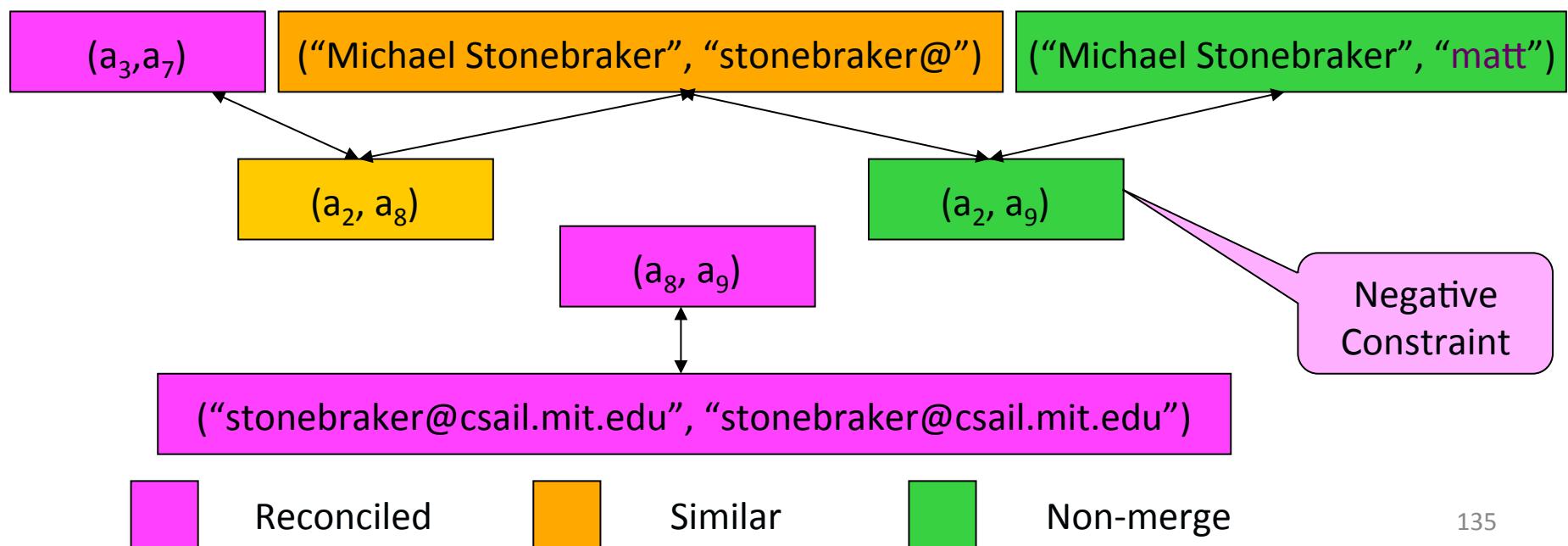
Enforce Constraints by Propagating Negative Information

- $a_2 = ("Michael\ Stonebraker",\ null,\ \{a_1,\ a_3\})$
 $a_3 = ("Eugene\ Wong",\ null,\ \{a_1,\ a_2\})$
 $a_7 = ("Eugene\ Wong",\ "eugene@berkeley.edu",\ \{a_8\})$
 $a_8 = (null,\ "stonebraker@csail.mit.edu",\ \{a_7\})$
 $a_9 = ("matt",\ "stonebraker@csail.mit.edu",\ null)$



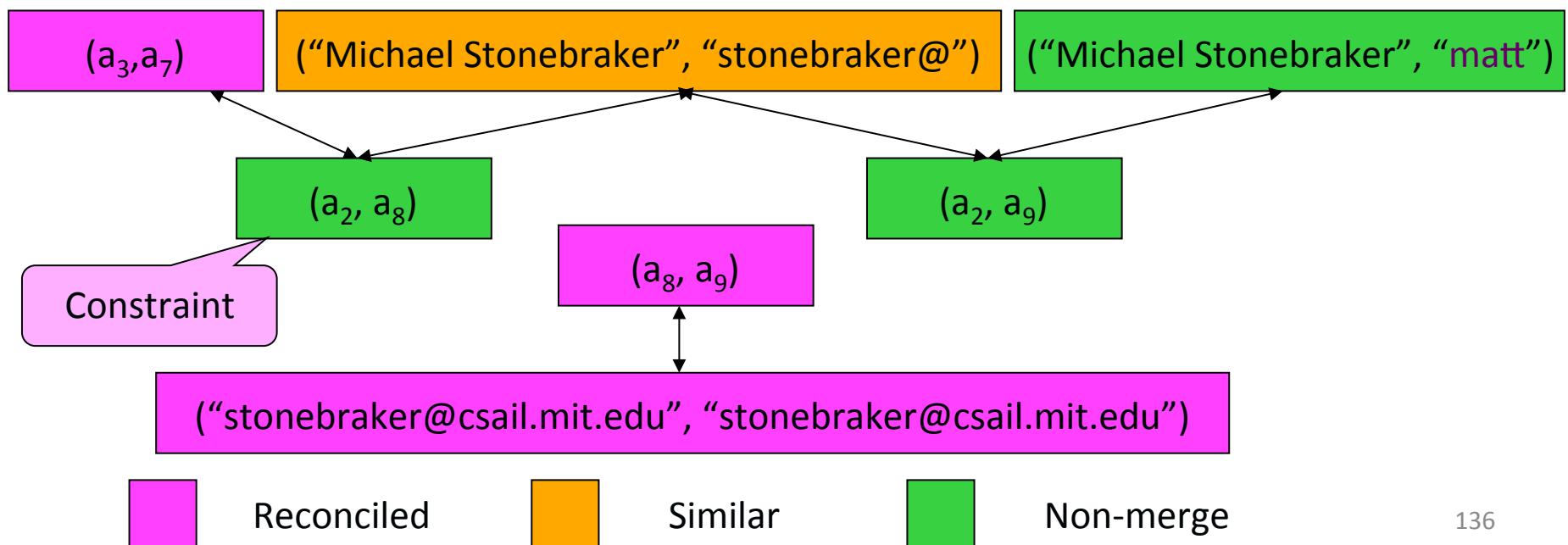
Enforce Constraints by Propagating Negative Information

- $a_2 = ("Michael\ Stonebraker", null, \{a_1, a_3\})$
 $a_3 = ("Eugene\ Wong", null, \{a_1, a_2\})$
 $a_7 = ("Eugene\ Wong", "eugene@berkeley.edu", \{a_8\})$
 $a_8 = (null, "stonebraker@csail.mit.edu", \{a_7\})$
 $a_9 = ("matt", "stonebraker@csail.mit.edu", null)$



Enforce Constraints by Propagating Negative Information

- $a_2 = ("Michael\ Stonebraker", null, \{a_1, a_3\})$
 $a_3 = ("Eugene\ Wong", null, \{a_1, a_2\})$
 $a_7 = ("Eugene\ Wong", "eugene@berkeley.edu", \{a_8\})$
 $a_8 = (null, "stonebraker@csail.mit.edu", \{a_7\})$
 $a_9 = ("matt", "stonebraker@csail.mit.edu", null)$

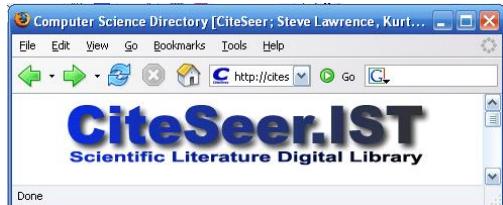


Collective Relational Clustering

[Bhattacharya & Getoor, TKDD07]

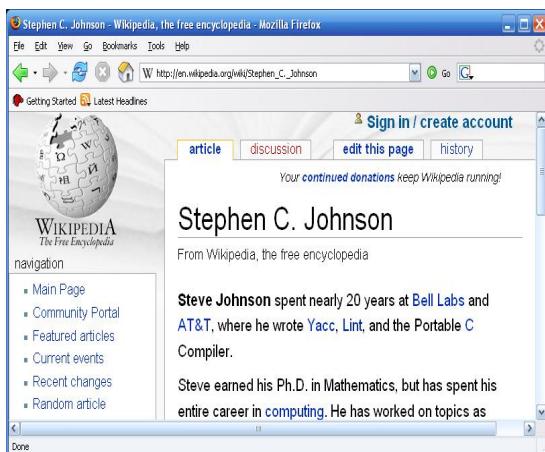
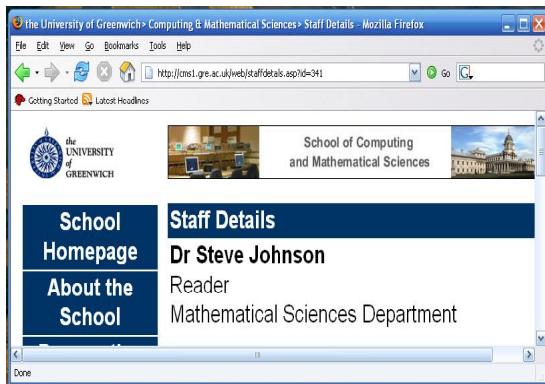
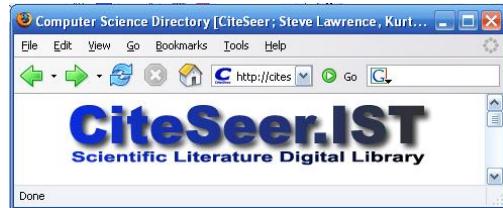
- Construct a graph where leaf nodes are individual mentions
- Perform hierarchical agglomerative clustering to merge clusters of mentions
- Similarity computed based on a combination of attribute and relational similarity
- When clusters are merged, update the similarities of any related clusters (clusters corresponding to mentions which co-occur with merged mentions)

A Motivating Example



- P1: "*JOSTLE: Partitioning of Unstructured Meshes for Massively Parallel Machines*", C. Walshaw, M. Cross, M. G. Everett, S. Johnson [J](#)
- P2: "*Partitioning Mapping of Unstructured Meshes to Parallel Machine Topologies*", C. Walshaw, M. Cross, M. G. Everett, S. Johnson, K. McManus [J](#)
- P3: "*Dynamic Mesh Partitioning: A Unified Optimisation and Load-Balancing Algorithm*", C. Walshaw, M. Cross, M. G. Everett
- P4: "*Code Generation for Machines with Multiregister Operations*", Alfred V. Aho, Stephen C. Johnson, Jefferey D. Ullman [J](#)
- P5: "*Deterministic Parsing of Ambiguous Grammars*", A. Aho, S. Johnson, J. Ullman [J](#)
- P6: "*Compilers: Principles, Techniques, and Tools*", A. Aho, R. Sethi, J. Ullman

A Motivating Example



P1: "*JOSTLE: Partitioning of Unstructured Meshes for Massively Parallel Machines*", C. Walshaw, M. Cross, M. G. Everett, S. Johnson

P2: "*Partitioning Mapping of Unstructured Meshes to Parallel Machine Topologies*", C. Walshaw, M. Cross, M. G. Everett, S. Johnson, K. McManus

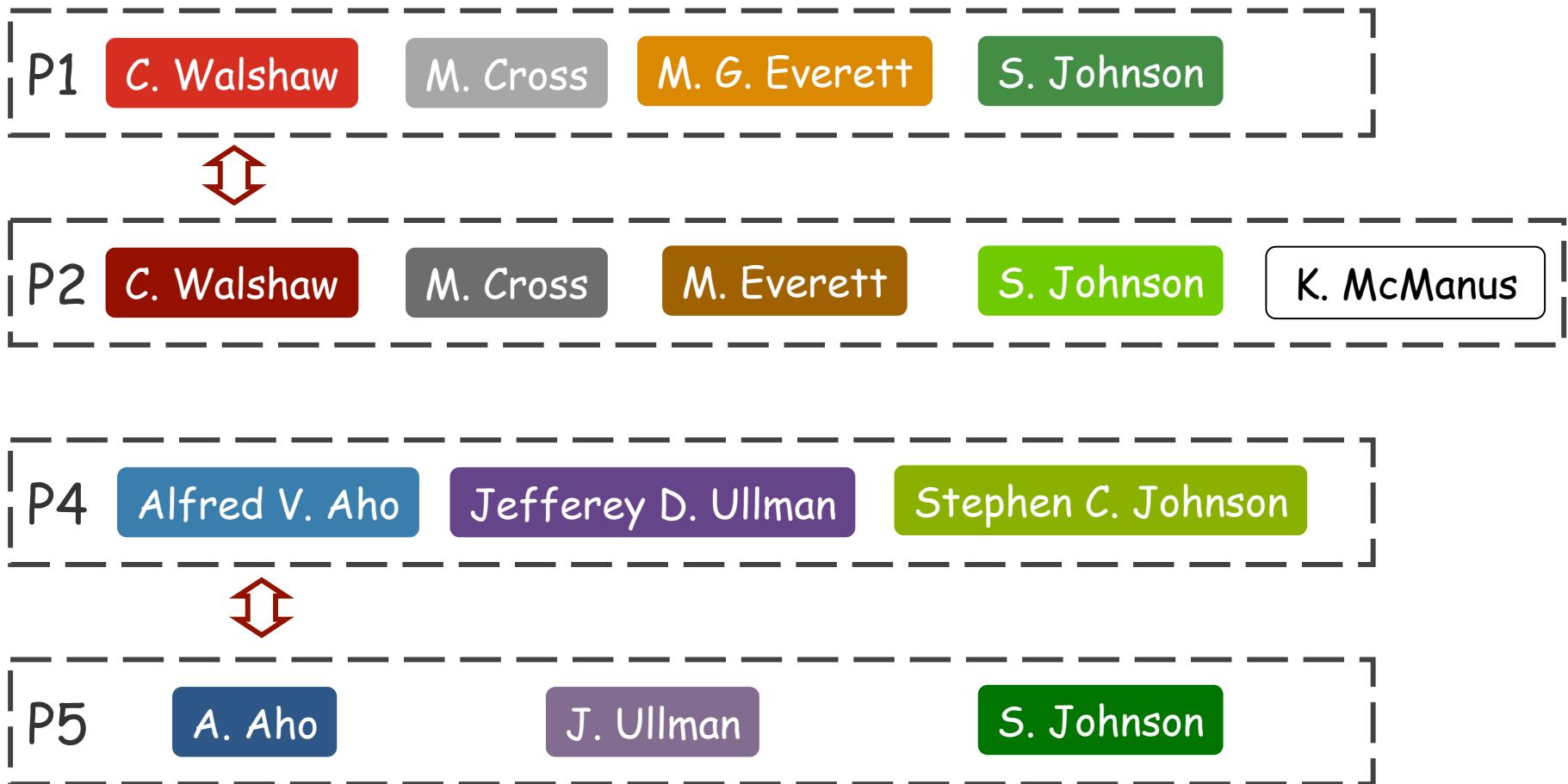
P3: "*Dynamic Mesh Partitioning: A Unified Optimisation and Load-Balancing Algorithm*", C. Walshaw, M. Cross, M. G. Everett

P4: "*Code Generation for Machines with Multiregister Operations*", Alfred V. Aho, Stephen C. Johnson, Jefferey D. Ullman

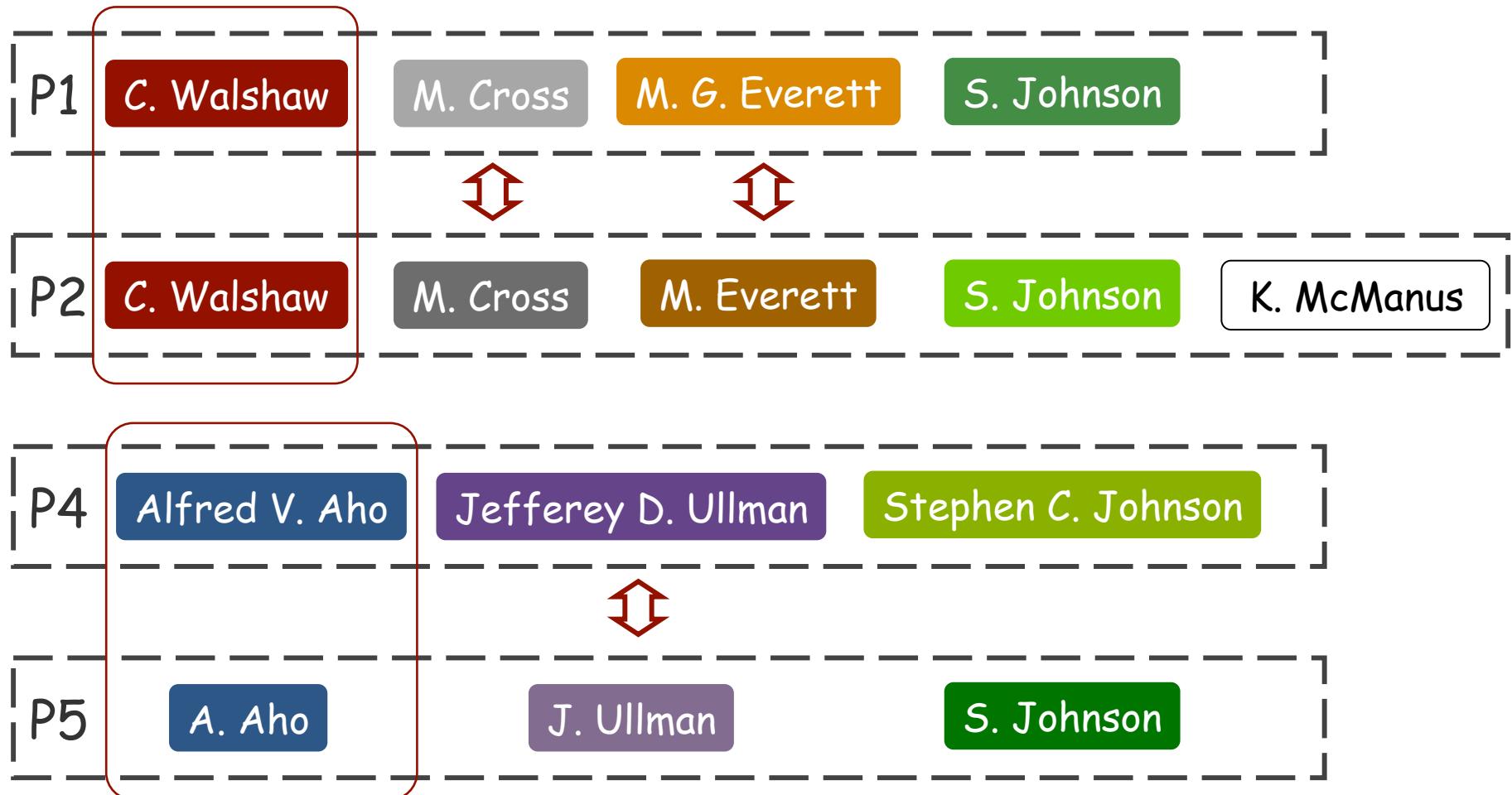
P5: "*Deterministic Parsing of Ambiguous Grammars*", A. Aho, S. Johnson, J. Ullman

P6: "*Compilers: Principles, Techniques, and Tools*", A. Aho, R. Sethi, J. Ullman

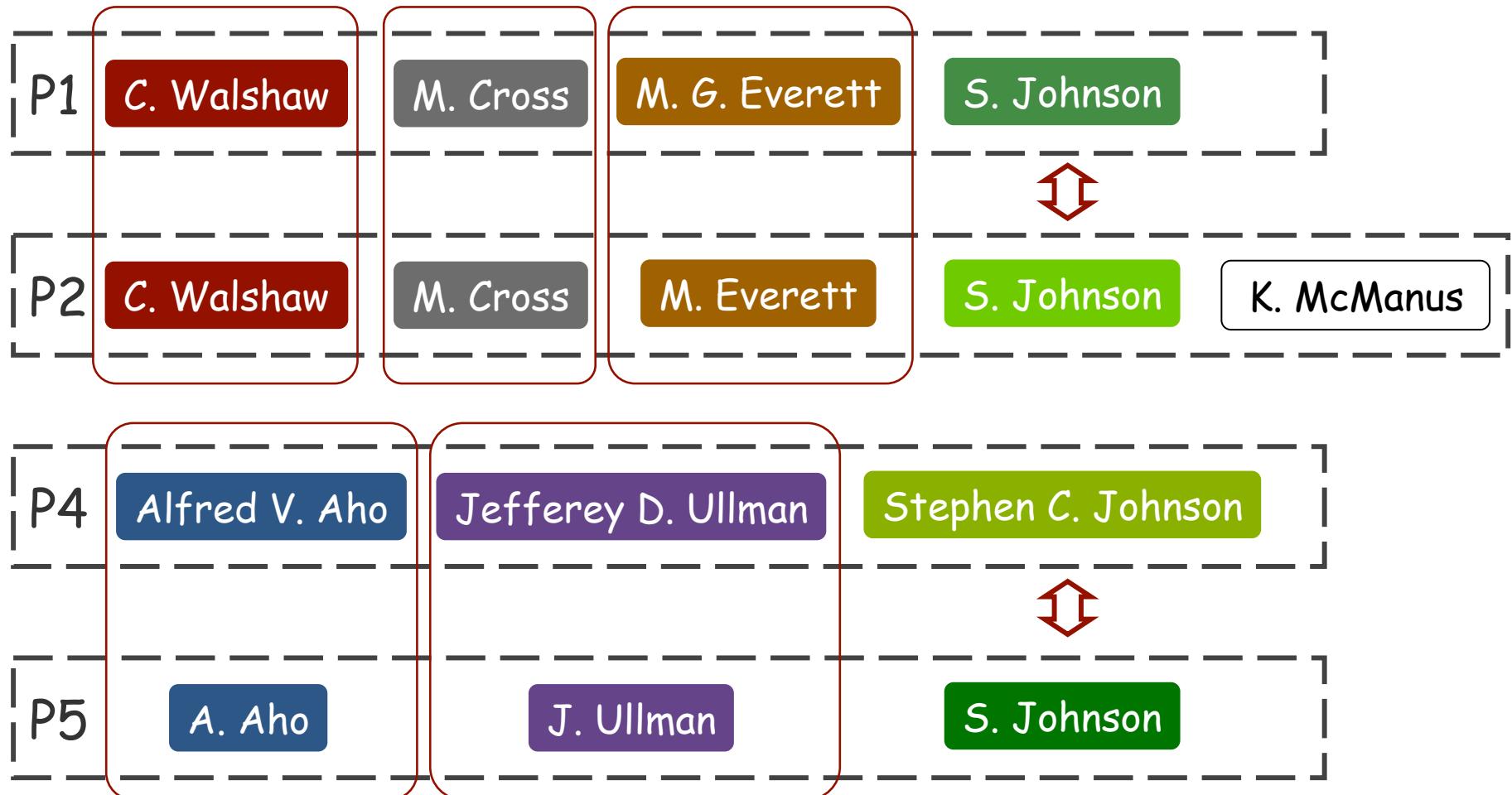
Relational Clustering for ER (RC-ER)



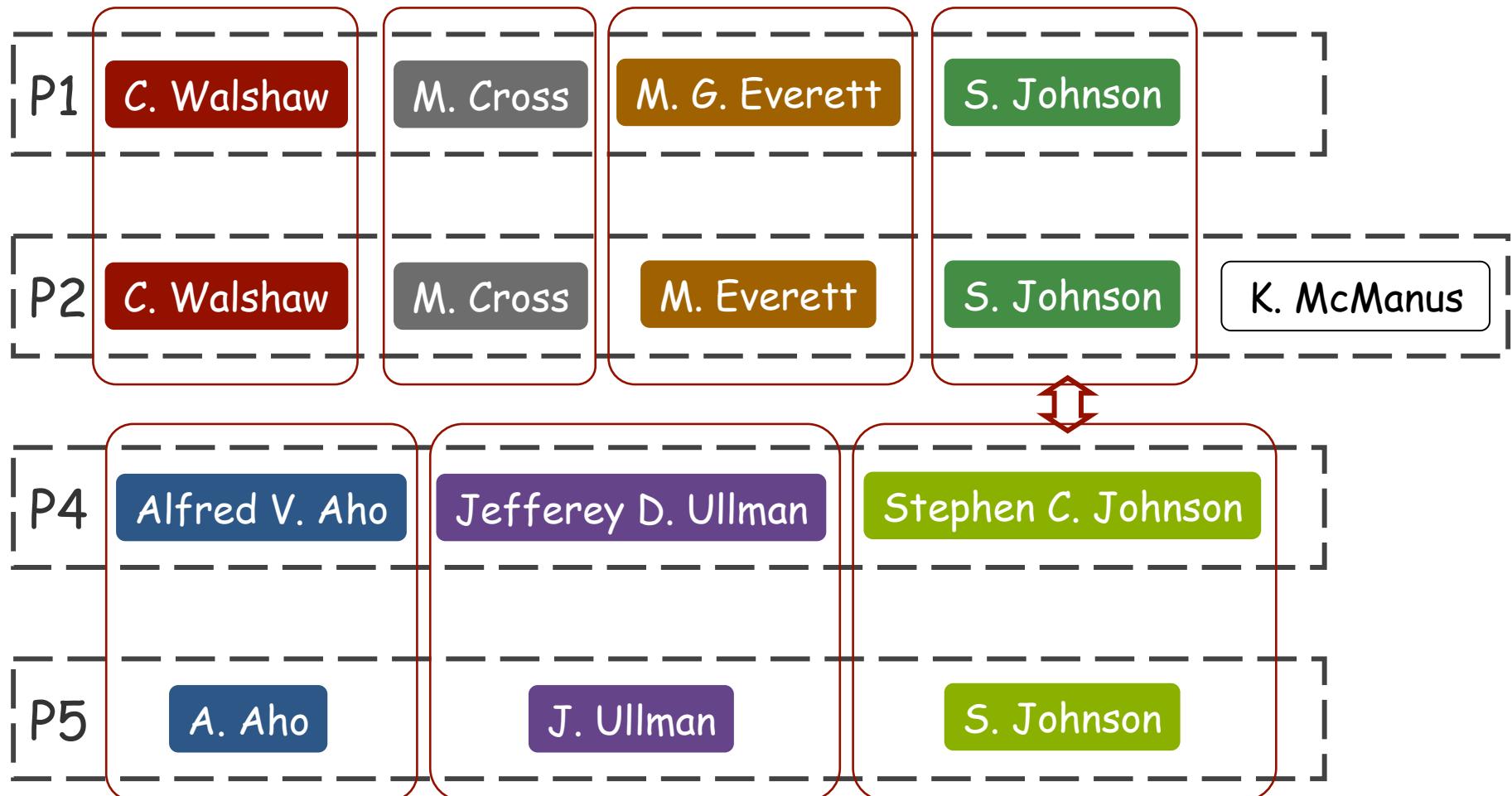
Relational Clustering for ER (RC-ER)



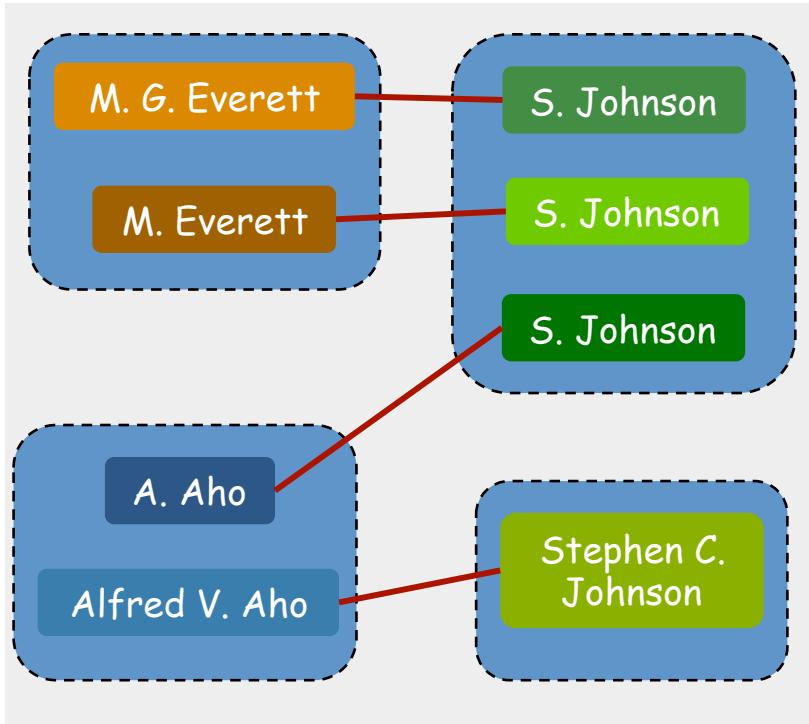
Relational Clustering for ER (RC-ER)



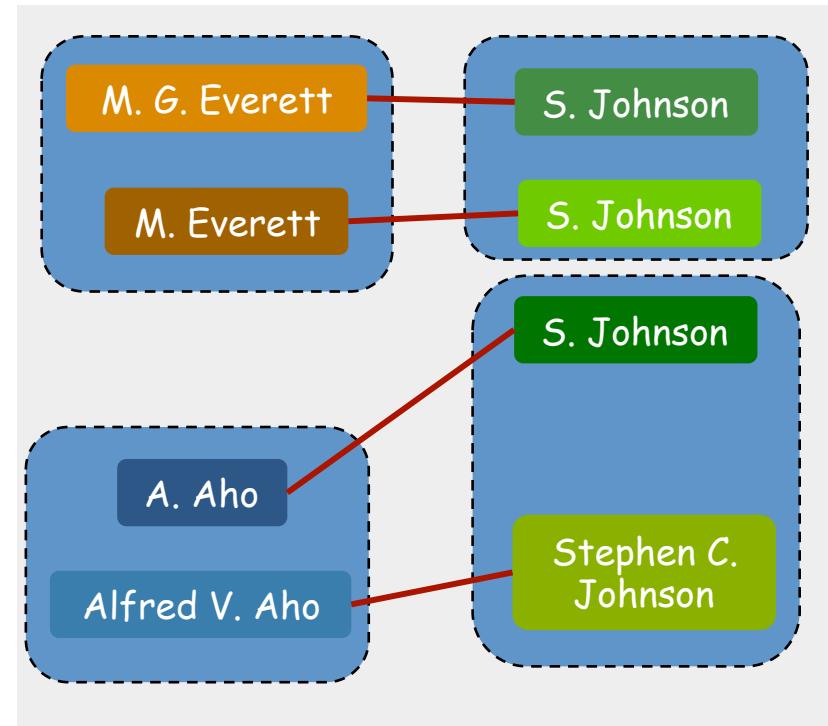
Relational Clustering for ER (RC-ER)



Cut-based Evaluation of Relational Clustering



Good separation of attributes
Many cluster-cluster relationships
➤ Aho-Johnson1, Aho-Johnson2,
Everett-Johnson1



Worse in terms of attributes
Fewer cluster-cluster relationships
➤ Aho-Johnson1, Everett-Johnson2

Objective Function

- Minimize:

$$\sum_i \sum_j w_A sim_A(c_i, c_j) + w_R sim_R(c_i, c_j)$$

weight for attributes similarity of attributes weight for relations Similarity based on relational edges between c_i and c_j

- Greedy clustering algorithm: merge cluster pair with max reduction in objective function

where for example

$$sim_A(c_i, c_j) = \sum_{a \in Attributes} sim(c_i^*, c_j^*) \quad \text{for cluster representative } c^*$$

and

$$sim_R(c_i, c_j) = sim_{jaccard}(N(c_i), N(c_j))$$

where $N(c)$ are the relational neighbors of c

Relational Clustering Algorithm

1. Find similar references using ‘blocking’
 2. Bootstrap clusters using attributes and relations
 3. Compute similarities for cluster pairs and insert into priority queue
 4. Repeat until priority queue is empty
 5. Find ‘closest’ cluster pair
 6. Stop if similarity below threshold
 7. If no negative constraints violated
 8. Merge to create new cluster
 9. Construct canonical cluster representative
 10. Update similarity for ‘related’ clusters
- $O(n k \log n)$ algorithm w/ efficient implementation

Similarity-propagation Approaches

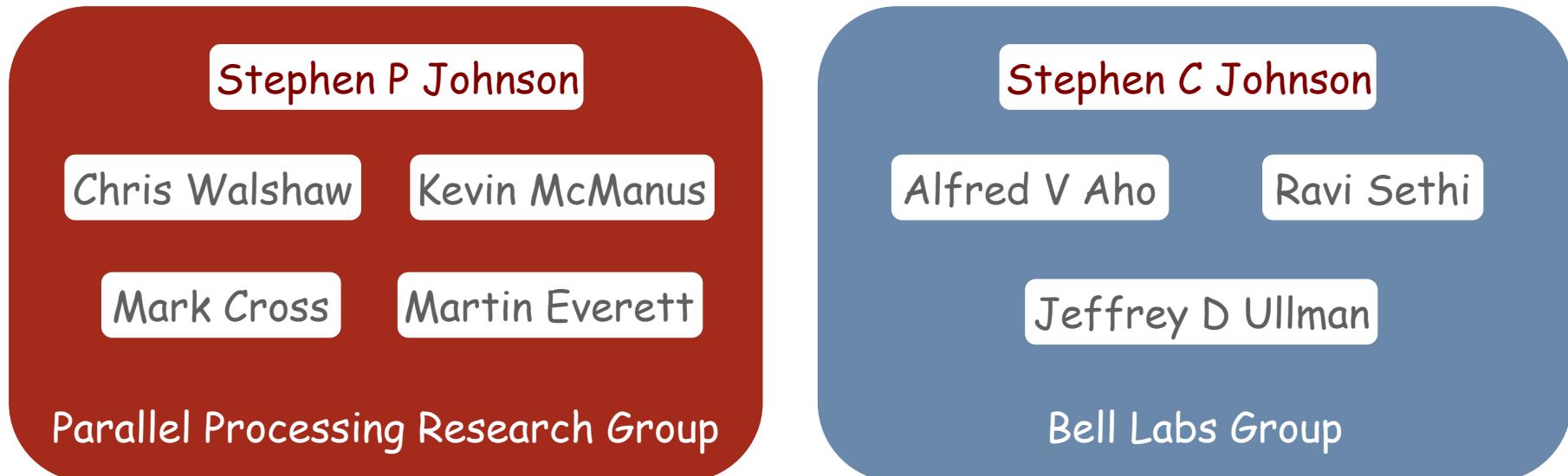
	Method	Notes	Constraints	Evaluation
RelDC [Kalashnikov et al, TODS06]	Reference disambiguation using using Relationship-based data cleaning (RelDC)	Model choice nodes identified using feature-based similarity	Context attraction measures the relational similarity	Accuracy and runtime for Author resolution and director resolution in Movie database
Reference Reconciliation [Dong et al, SIGMOD05]	Dependency Graph for propagating similarities + enforce non-match constraints	Reference enrichment Explicitly handle missing values Parameters set by hand	Both positive and negative constraints	Precision/Recall, F1 on personal information management data (PIM), Cora dataset
Collective Relational Clustering [Bhattacharya & Getoor, TKDD07]	Modified hierarchical agglomerative clustering approach	Constructs canonical entity as merges are made	Focus on coauthor resolution and propagation	Precision/Recall, F1 on three bibliographic datasets: CiteSeer, ArXiv, and BioBase, and synthetic data ¹⁴⁷

PROBABILISTIC MODELS: GENERATIVE APPROACHES

Generative Probabilistic Approaches

- Probabilistic semantics based on Directed Models
 - Model dependencies between match decisions in a generative manner
 - Disadvantage: acyclicity requirement
- Variety of approaches
 - Based on Latent Dirichlet Allocation, Bayesian Networks
- Examples
 - Latent Dirichlet Allocation [Bhattacharya & Getoor, SDM07]
 - Probabilistic Relational Models [Pasula et al, NIPS02]

LDA for Entity Resolution: Discovering Groups from Co-Occurrence Relations



P1: C. Walshaw, M. Cross, M. G. Everett,
S. Johnson

P2: C. Walshaw, M. Cross, M. G. Everett,
S. Johnson, K. McManus

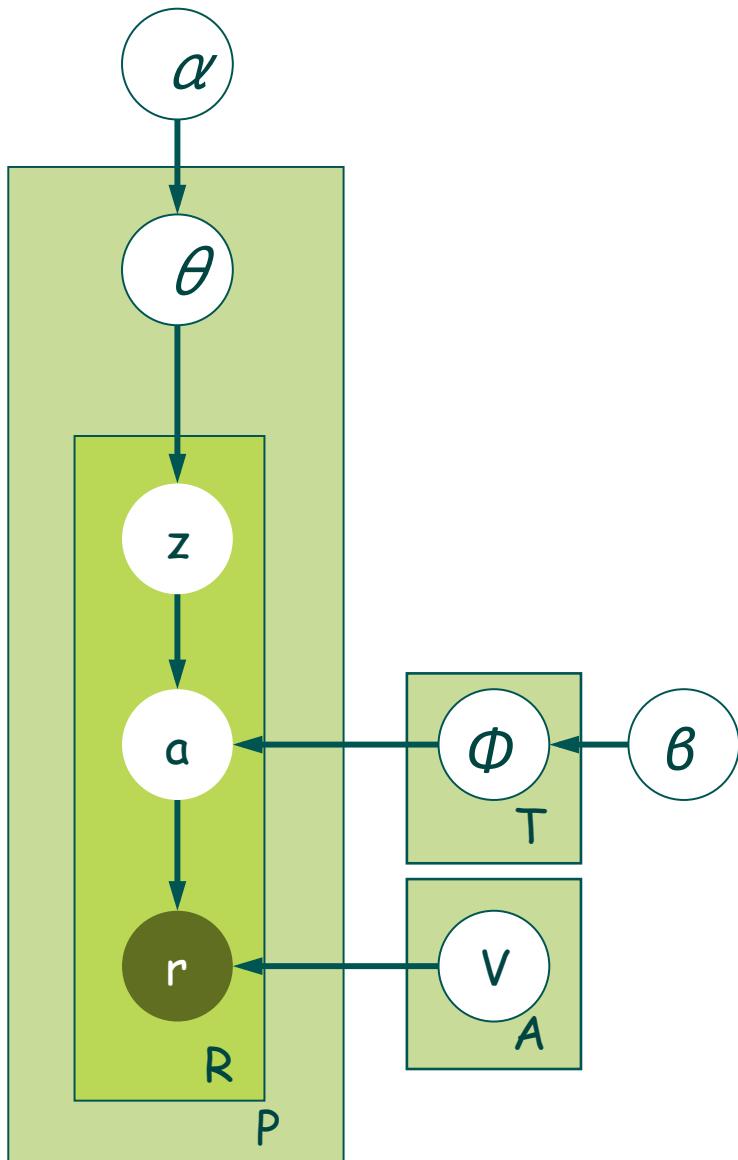
P3: C. Walshaw, M. Cross, M. G. Everett

P4: Alfred V. Aho, **Stephen C. Johnson**,
Jefferey D. Ullman

P5: A. Aho, **S. Johnson**, J. Ullman

P6: A. Aho, R. Sethi, J. Ullman

LDA-ER Model



- Entity label a and group label z for each reference r
- Θ : 'mixture' of groups for each co-occurrence
- ϕ_z : multinomial for choosing entity a for each group z
- V_a : multinomial for choosing reference r from entity a
- Dirichlet priors with α and β

Inference using blocked Gibbs sampling for efficiency (and improved accuracy)

Generative Approaches

	Method	Learning/Inference Method	Evaluation
[Li, Morie, & Roth, AAAI 04]	Generative model for mentions in documents	Truncated EM to learn parameters and MAP inference for entities (unsupervised)	F1 on person names, locations and organizations in TREC dataset
Probabilistic Relational Models [Pasula et al., NIPS03]	Probabilistic Relational Models	Parameters learned on separated corpora, inference done using MCMC	% of correctly identified clusters on subsets of CiteSeer data
Latent Dirichlet Allocation [Bhattacharya & Getoor, SDM06]	Latent-Dirichlet Allocation Model	Blocked Gibbs Sampling Unsupervised approach	Precision/ Recall/F1 on CiteSeer and HEP data

PROBABILISTIC MODELS: UNDIRECTED APPROACHES

Undirected Probabilistic Approaches

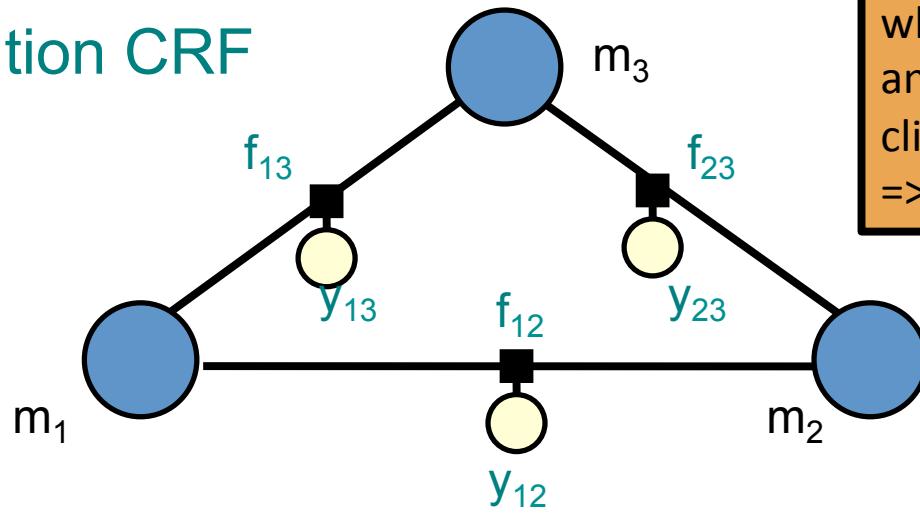
- Probabilistic semantics based on Markov Networks
 - Advantage: no acyclicity requirements
- In some cases, syntax based on first-order logic
 - Advantage: declarative
- Examples
 - **Conditional Random Fields (CRFs)** [McCallum & Wellner, NIPS04]
[Sarawagi & Cohen, NIPS04]
 - Markov Logic Networks (MLNs) [Singla & Domingos, ICDM06]
 - Probabilistic Soft Logic [Broeckeler & Getoor, UAI10]

[Lafferty, McCallum, Pereira, ICML01]

[McCallum & Wellner, NIPS04]

- x_1 is a mention, J. Smith
- x_2 is a mention, John
- x_3 is a mention, John Smith
- y_{12}, y_{23}, y_{13} are binary variables representing match or not match
- f_{12} is a factor between x_1, x_2, y_{12}
- f_{23} is a factor between x_2, x_3, y_{23}
- f_{13} is a factor between x_1, x_3, y_{13} ,

Entity Resolution CRF



inference in CRF is equivalent to graph partitioning in graph where nodes are mentions and edges weights are log clique potentials over nodes
=> correlation clustering!

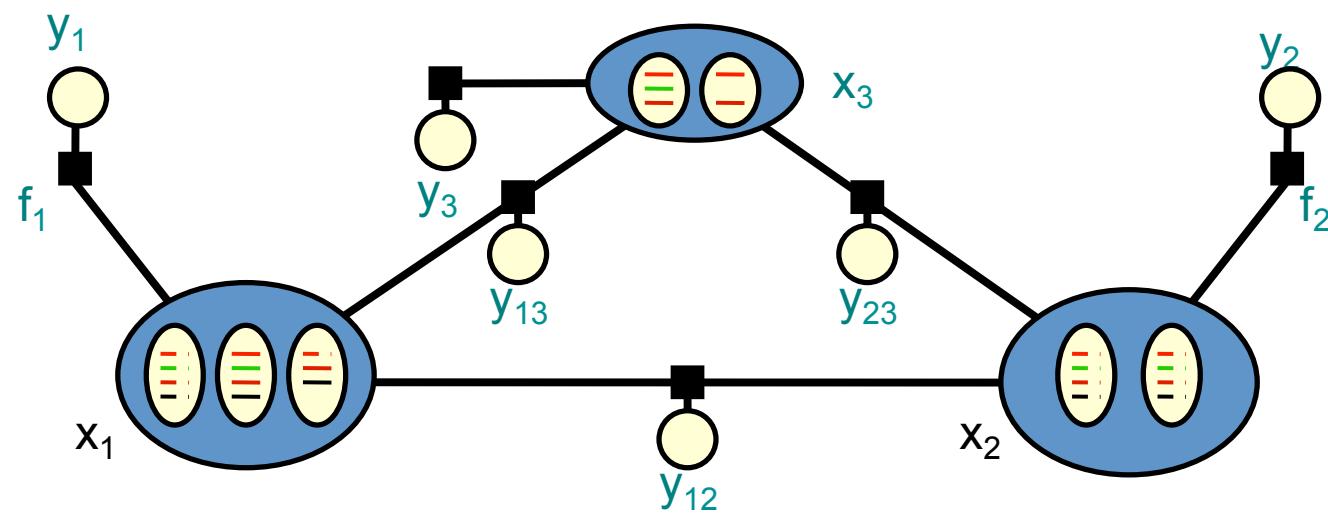
$$P(y|x) = \frac{1}{Z_x} \prod_f f(x \in f, y \in f)$$

[Based on slides from Andrew McCallum]

[Wick et al., SDM09]

- x_1 is a set of mentions {J. Smith, John, John Smith}
- x_2 is a set of mentions {Amanda, A. Jones}
- f_{12} is a factor between x_1/x_2
- y_{12} is a binary variable indicating a match (no)
- f_1 is a factor over cluster x_1
- y_1 is a binary variable indicating match (yes)
- Entity/attribute factors omitted for clarity

Approximate inference using
agglomerative clustering



Entity Reference and Canonicalization

Undirected Probabilistic Approaches

- Probabilistic semantics based on Markov Networks
 - Advantage: no acyclicity requirements
- In some cases, syntax based on first-order logic
 - Advantage: declarative
- Examples
 - Conditional Random Fields (CRFs) [McCallum & Wellner, NIPS04]
[Sarawagi & Cohen, NIPS04]
 - **Markov Logic Networks (MLNs)** [Singla & Domingos, ICDM06]
 - Probabilistic Soft Logic [Broeckeler & Getoor, UAI10]

Markov Logic

- A logical KB is a set of **hard constraints** on the set of possible worlds
- Make them **soft constraints**; when a world violates a formula, it becomes less probable but not impossible
- Give each formula a **weight**
 - Higher weight \Rightarrow Stronger constraint

$$P(world) \propto \exp\left(\sum \text{weights of formulas it satisfies}\right)$$

[Richardson & Domingos, 06]
163

Markov Logic

- A **Markov Logic Network (MLN)** is a collection of first-order formulas **F** and real-valued weights **w**

$$P(X) = \frac{1}{Z} \exp\left(\sum_{i \in F} w_i n_i(x) \right)$$

Diagram annotations:

- A red arrow points from the text "Normalization Constant" to the variable Z .
- A blue arrow points from the text "Iterate over all first-order MLN formulas" to the summand $w_i n_i(x)$.
- A blue box contains the text "# true groundings of *i*th clause" with a blue arrow pointing to the term $n_i(x)$.

[Richardson & Domingos, 06]
164

ER Problem Formulation in MLNs

- **Given**

- A DB of records representing mentions of entities in the real world, e.g. paper mentions
- A set of fields e.g. author, title, venue
- Each record represented as a set of typed predicates e.g. $\text{HasAuthor}(\text{paper}, \text{author})$, $\text{HasVenue}(\text{paper}, \text{venue})$

- **Goal**

- Determine which of the records/fields refer to the same underlying entity

Handling Equality

- Introduce *Equals(x,y)* or $x = y$
- Introduce the axioms of equality
 - Reflexivity: $x = x$
 - Symmetry: $x = y \Rightarrow y = x$
 - Transitivity: $x = y \wedge y = z \Rightarrow z = x$
 - Predicate Equivalence:

$$x_1 = x_2 \wedge y_1 = y_2 \Rightarrow (R(x_1, y_1) \Leftrightarrow R(x_2, y_2))$$

Positive, Soft Evidence

- Introduce **reverse predicate equivalence**
- Same relation with the same entity gives evidence about two entities being same

$$R(x_1, y_1) \wedge R(x_2, y_2) \wedge x_1 = x_2 \Rightarrow y_1 = y_2$$

- Not true logically, but gives useful information
- Example

$$\text{HasAuthor}(P1, J. \text{ Cox}) \wedge \text{HasAuthor}(P2, \text{Cox J.}) \wedge P1 = P2 \Rightarrow (J. \text{ Cox} = \text{Cox J.})$$

Entity Resolution

- Field similarities measured by presence/absence of words in common

$$\text{HasWord}(f_1, w_1) \wedge \text{HasWord}(f_2, w_2) \wedge \text{HasField}(r_1, f_1) \wedge \\ \text{HasField}(r_2, f_2) \wedge w_1 = w_2 \Rightarrow r_1 = r_2$$

- Example

$$\text{HasWord}(J.~Cox, Cox) \wedge \text{HasWord}(Cox J., Cox) \wedge \text{HasAuthor}(P1, \\ J.~Cox) \wedge \text{HasAuthor}(P2, Cox J.) \wedge (Cox = Cox) \Rightarrow (P1 = P2)$$

- Additional Constraints

$$\text{HasAuthor}(c, a_1) \wedge \text{HasAuthor}(c, a_2) \Rightarrow \text{Coauthor}(a_1, a_2) \\ \text{Coauthor}(a_1, a_2) \wedge \text{Coauthor}(a_3, a_4) \wedge a_1 = a_3 \Rightarrow a_2 = a_4$$

Inference

- Use cheap heuristics (e.g. TFIDF based similarity) to identify plausible pairs
- Inference/learning over plausible pairs
- Inference method: lazy grounding + MaxWalkSAT
- Learning: supervised and transfer (learn/hand set on one domain and transferred to another domain)

Undirected Probabilistic Approaches

- Probabilistic semantics based on Markov Networks
 - Advantage: no acyclicity requirements
- In some cases, syntax based on first-order logic
 - Advantage: declarative
- Examples
 - Conditional Random Fields (CRFs) [McCallum & Wellner, NIPS04]
[Sarawagi & Cohen, NIPS04]
 - Markov Logic Networks (MLNs) [Singla & Domingos, ICDM06]
 - **Probabilistic Soft Logic** [Broecheler & Getoor, UAI10]

Probabilistic Soft Logic

- Declarative language for defining **hinge-loss Markov random field** (CCMRF) using first-order logic (FOL)
- Soft logic: truth values in $[0,1]$
- Logical operators relaxed using Lukasiewicz t-norms
- Mechanisms for incorporating similarity functions, and reasoning about sets
- MAP inference is a **convex optimization**
- Efficient sampling method for marginal inference

FOL to HL-MRF

- PSL converts a weighted rule into potential functions by penalizing its **distance to satisfaction**, $d(g, x)$
- The distribution over truth values is

$$\Pr(x) = \frac{1}{Z} \exp \left(-\sum_{r \in P} \sum_{g \in G(r)} w_r d(g, x) \right)$$

w_r : weight of rule r

$G(r)$: all groundings of rule r

P : PSL program

Summary of PSL for ER

- Declarative representation
- Undirected graphical model
- *Continuous, convex* optimization
- Significant scaling improvements [Bach et al, NIPS12, Bach et al, UAI13]

<http://psl.umiacs.umd.edu>

Undirected Approaches

	Method	Learning/Inference Method	Evaluation
[McCallum & Wellner, NIPS04]	Conditional Random Fields (CRFs) capturing transitivity constraints	Graph partitioning (Boykov et al. 1999), performed via correlation clustering	F1 on DARPA MUC & ACE datasets
[Singla & Domingos, ICDM06]	Markov Logic Networks (MLNs)	Supervised learning and inference using MaxWalkSAT & MCMC	Conditional Log-likelihood and AUC on Cora and BibServ data
[Broeckeler & Getoor, UAI10]	Probabilistic Soft Logic (PSL)	Supervised learning and inference using continuous optimization	Precision/Recall/F1 Ontology Alignment

HYBRID APPROACHES

Hybrid Approaches

- Constraint-based approaches explicitly encode relational constraints
 - They can be formulated as hybrid of constraints and probabilistic models
- Examples
 - Constraint-based Entity Matching [Shen, Li & Doan, AAAI05]
 - Performs relaxation labeling
 - Dedupalog [Arasu, Re, Suciu, ICDE09]
 - Combines correlation clustering with declarative representation

Dedupalog [Arasu et al., ICDE09]

PaperRef(id, title, conference, publisher, year)
Wrote(id, authorName, Position)

Data to be
deduplicated

TitleSimilar(title1,title2)
AuthorSimilar(author1,author2)

Thresholded Similarity

Step (0) Create initial approximate matches; this is input to Dedupalog.

Step (1) Declare the entities “*Cluster Papers, Publishers, & Authors*”

```
Paper!(id)    :- PaperRef(id,-,-,-)
Publisher!(p) :- PaperRef(-,-,-,p,-)
Author!(a)   :- Wrote(-,a,-)
```

Step (2) Declare Clusters

Input in the DB

PaperRef(id, title, conference, publisher, year)
Wrote(id, authorName, Position)

“Cluster papers, publishers, and authors”

TitleSimilar(title1,title2)
AuthorSimilar(author1,author2)

Paper!(id) :- PaperRef(id,-,-,-)
Publisher!(p) :- PaperRef(-,-,-,p,-)
Author!(a) :- Wrote(-,a,-)

Clusters are *declared* using * (like IDBs or Views): These are output

Author*(a₁,a₂) <-> AuthorSimilar(a₁,a₂)

“Cluster authors with similar names”

*IDBs are equivalence relations:

Symmetric, Reflexive , & Transitively-Closed Relations: i.e., *Clusters*

Inference in Dedupalog via CC

Semantics: Translate a Dedupalog Program to a set of graphs

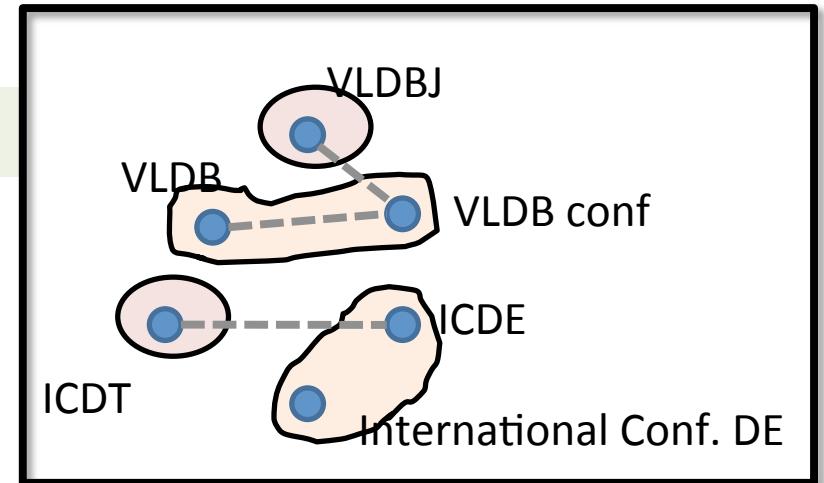
Nodes are references (in the ! Relation)

Entity References: Conference!(c)

Conference^{*}(c₁, c₂) <-> ConfSim(c₁, c₂)

— Positive edges

[−] Negative edges are implicit

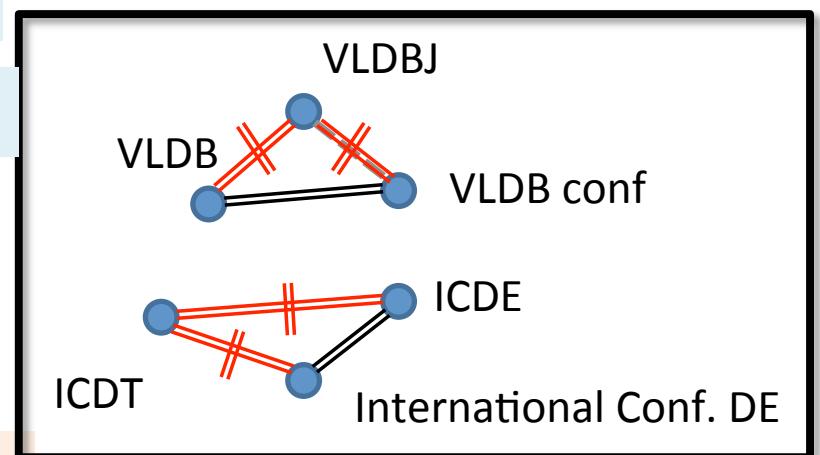
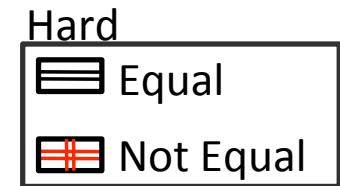
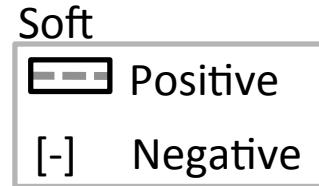


Correlation Clustering

Conference*(c_1, c_2) \leftarrow ConfSim(c_1, c_2)

Conference*(c_1, c_2) \leq ConfEQ(c_1, c_2)

\neg **Conference***(c_1, c_2) \leq ConfNEQ(c_1, c_2)



1. Pick a random order of edges
2. **While** there is a soft edge **do**
 1. Pick first soft edge in order
 2. If [---] turn into [==]
 3. Else is [-] turn into [!=]
 4. Deduce labels
3. **Return** Transitively closed subsets

Hybrid Approaches

	Method	Evaluation
Constraint-based Entity Matching [Shen, Li & Doan, AAAI05]; builds on (Li, Morie, & Roth, AI Mag 2004)	<p>Two layer model:</p> <p>Layer 1: Generative model for data sets that satisfy constraints;</p> <p>Layer 2: EM algorithm and the relaxation labeling algorithm to perform matching. In each iteration, use EM to estimate parameters of the generative model and a matching assignment, then employs relaxation labeling to exploit the constraints</p>	Researchers and IMDB with noise added
Dedupalog [Arasu, Re, Suciu, ICDE09]	Declarative specification for rich collection of constraints with nice syntactic sugar added to datalog for ER. Inference: Correlation clustering+ voting	Precision/Recall on Cora, subset of ACM dataset

Summary: Collective Approaches

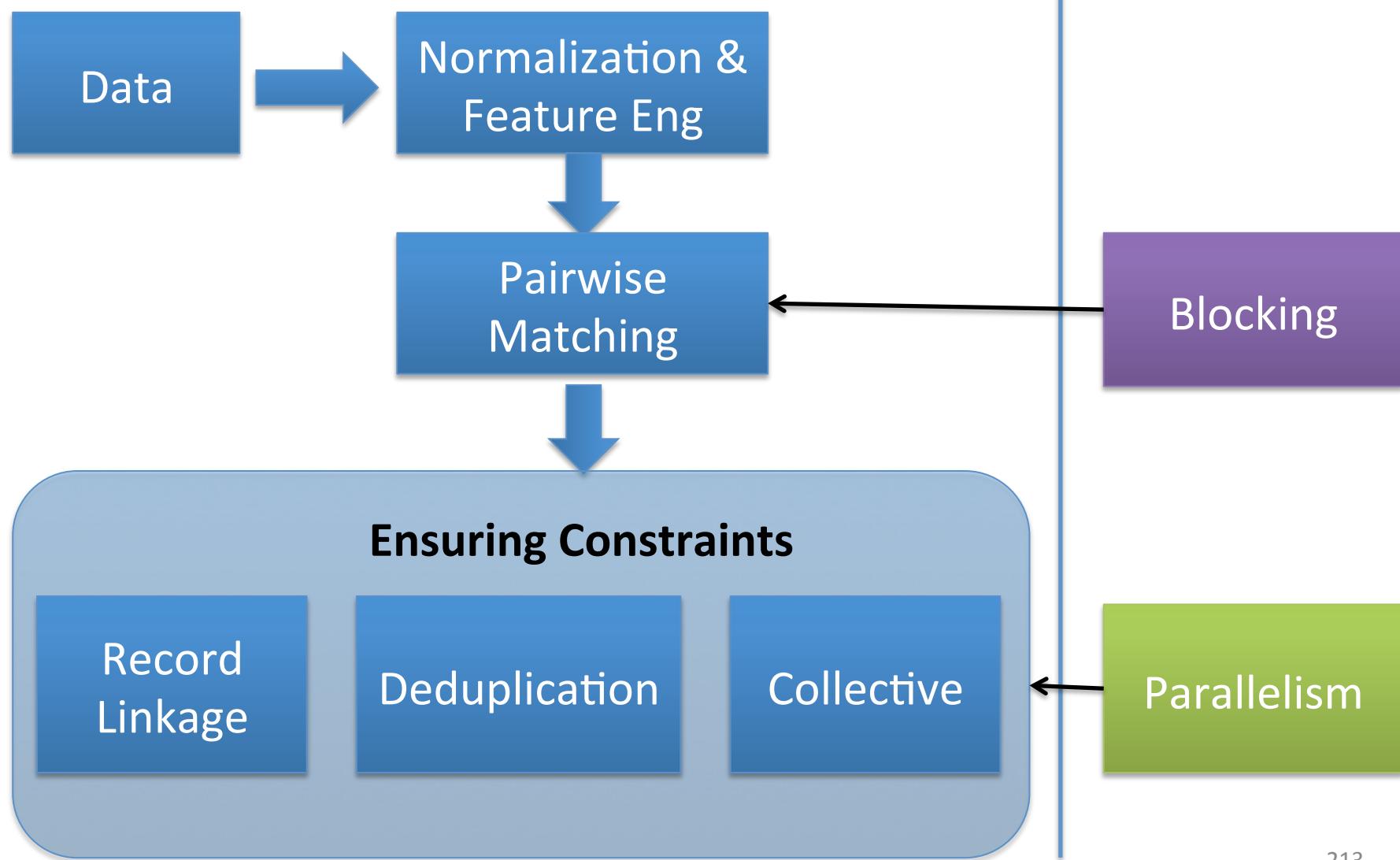
- Decisions for cluster-membership depends on other clusters
 - Similarity propagation
 - Agglomerative clustering
 - Probabilistic Models
 - Generative Models
 - Undirected Models
 - Hybrid Approaches
- Non-probabilistic approaches often scale better than probabilistic approaches
- Undirected/constraint-based models are often easier to specify
- Scaling active area of research

PART 3

SCALING ER TO BIG DATA

Part 2: Algorithmic Foundations

Part 3: Scaling

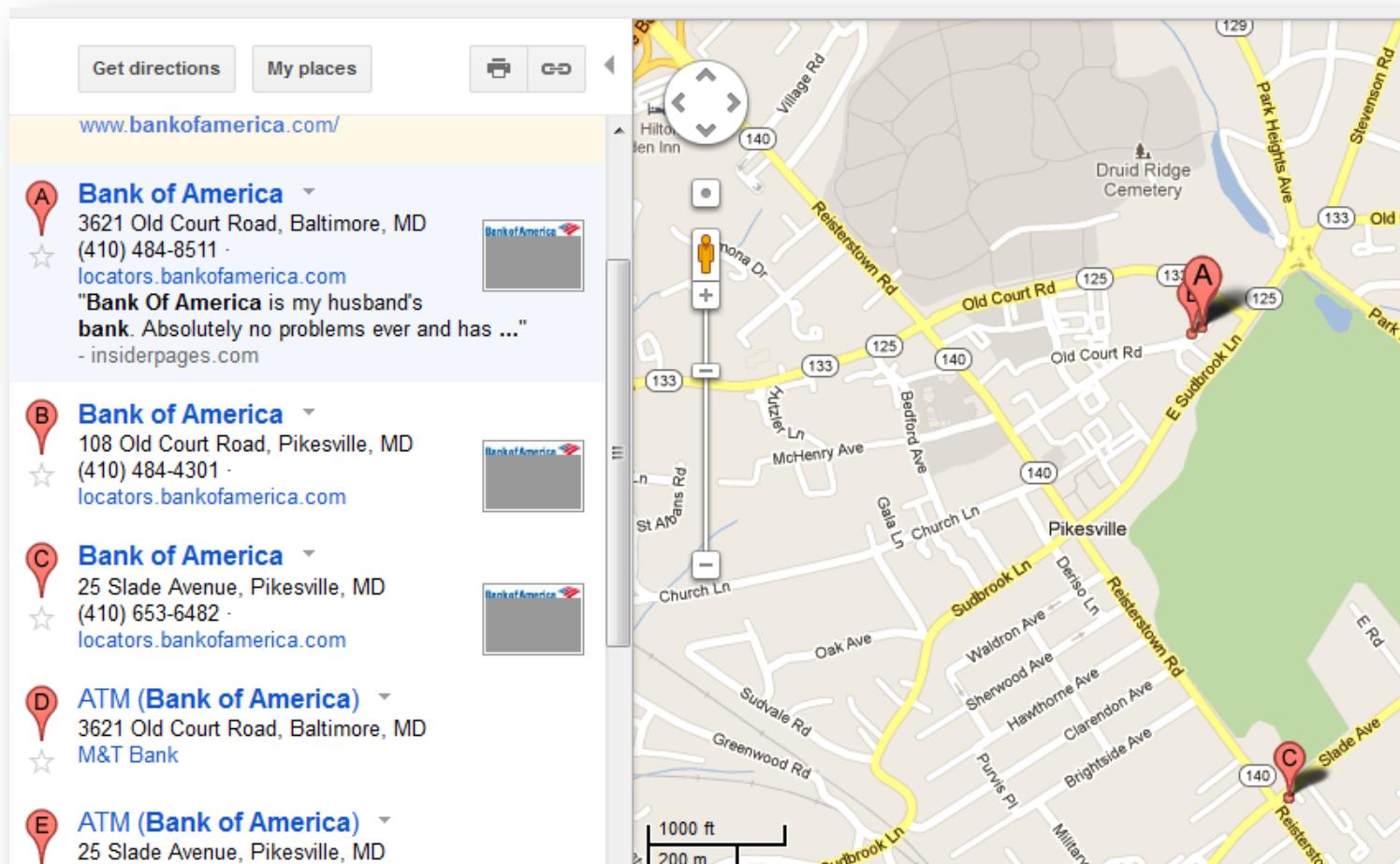


Blocking: Motivation

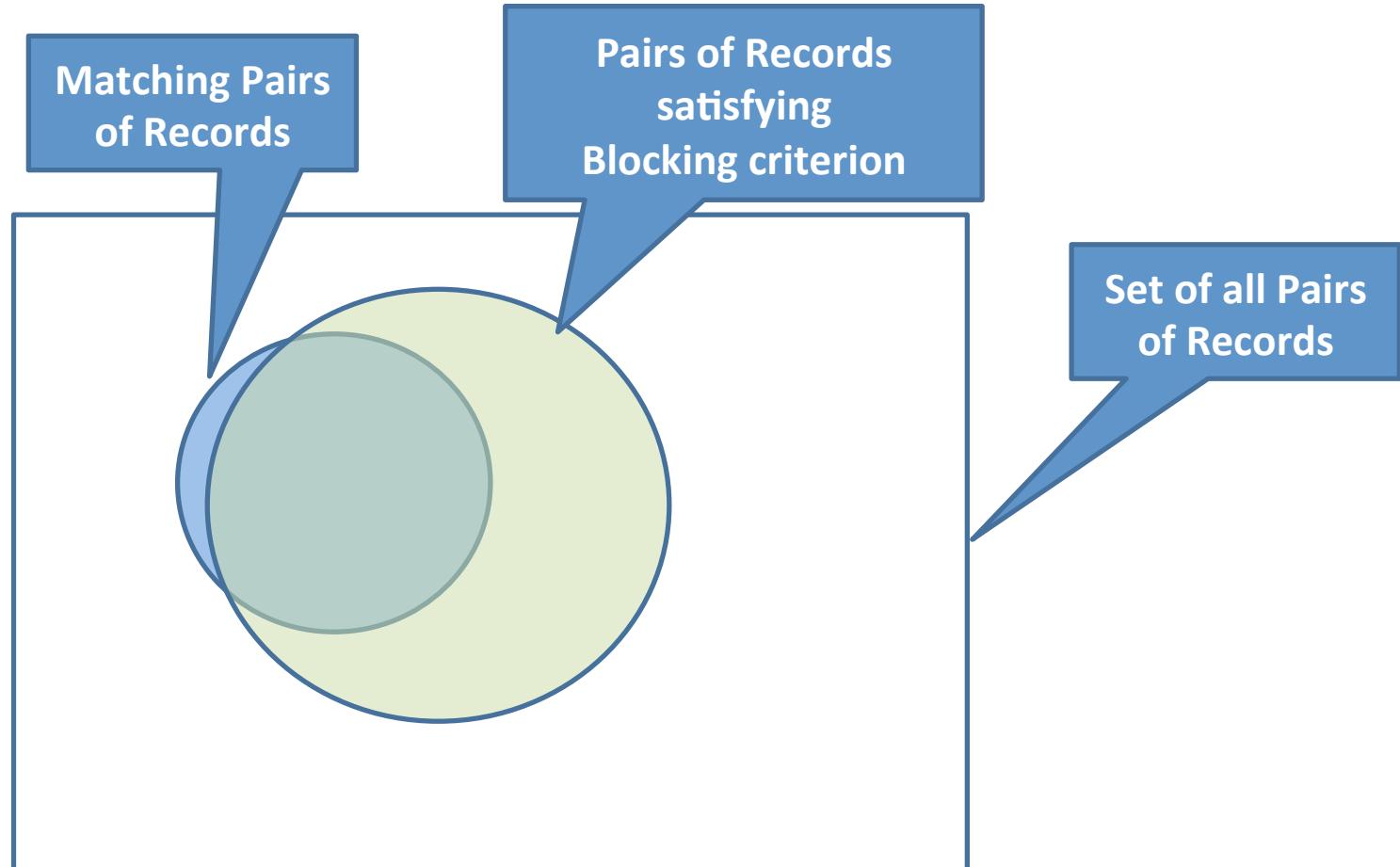
- Naïve pairwise: $|R|^2$ pairwise comparisons
 - 1000 business listings each from 1,000 different cities across the world
 - 1 trillion comparisons
 - 11.6 days (if each comparison is 1 μs)
- Mentions from different cities are unlikely to be matches
 - **Blocking Criterion: City**
 - 1 billion comparisons
 - 16 minutes (if each comparison is 1 μs)

Blocking: Motivation

- Mentions from different cities are unlikely to be matches
 - May miss potential matches



Blocking: Motivation



Blocking Algorithms 1

- Hash based blocking
 - Each block C_i is associated with a hash key h_i .
 - Mention x is hashed to C_i if $\text{hash}(x) = h_i$.
 - Within a block, all pairs are compared.
 - Each hash function results in disjoint blocks.
- What *hash* function?
 - Deterministic function of attribute values
 - Boolean Functions over attribute values
[Bilenko et al ICDM'06, Michelson et al AAAI'06,
Das Sarma et al CIKM '12]
 - **minHash** (min-wise independent permutations)
[Broder et al STOC'98]

Blocking Algorithms 2

- Pairwise Similarity/Neighborhood based blocking
 - Nearby nodes according to a similarity metric are clustered together
 - Results in non-disjoint canopies.
- Techniques
 - Sorted Neighborhood Approach [Hernandez et al SIGMOD'95]
 - Canopy Clustering [McCallum et al KDD'00]

Blocking Techniques

- **Predicate blocking**
- Locality Sensitive Hashing
- Canopy Clustering

Simple Blocking: Inverted Index on a Predicate

Examples of blocking predicates (keys):

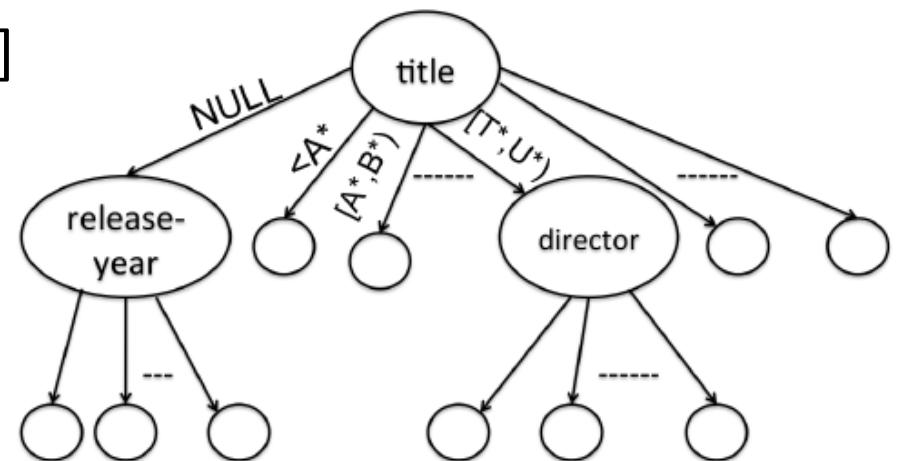
- First three characters of last name
- City + State + Zip
- Character or Token n-grams
- Minimum infrequent n-grams

Learning Optimal Blocking Functions

- Using one or more blocking predicates may be insufficient
 - 2,376,206 American's shared the surname Smith in the 2000 US
 - NULL values may create large blocks.
- Solution: Construct blocking predicates by combining simple predicates

Complex Blocking Predicates

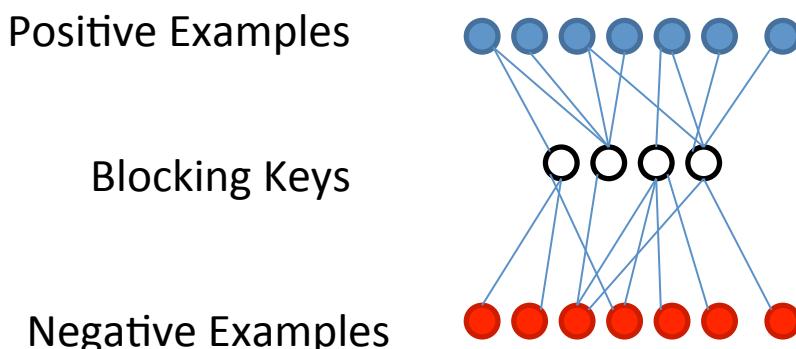
- Conjunction of predicates [Michelson et al AAAI'06, Bilenko et al ICDM'06]
 - {City} AND {last four digits of phone}
- Chain-trees [Das Sarma et al CIKM'12]
 - If ($\{\text{City}\}$ = NULL or LA) then {last four digits of phone} AND {area code}
else {last four digits of phone} AND {City}
- BlkTrees [Das Sarma et al CIKM'12]



Learning an Optimal predicate

[Bilenko et al ICDM '06]

- Find k blocking predicates that eliminate the most non-matches, while retaining almost all matches.
 - Need a training set of positive and negative pairs
- Algorithm Idea: Red-Blue Set Cover



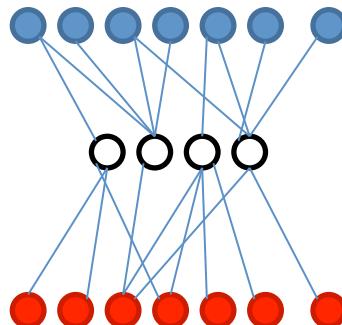
Pick k Blocking keys such that

- (a) At most ϵ blue nodes are not covered
- (b) Number of red nodes covered is minimized

Learning an Optimal function [Bilenko et al ICDM '06]

- Algorithm Idea: Red-Blue Set Cover

Positive Examples



Blocking Keys

Negative Examples

Pick k Blocking keys such that

- (a) At most ϵ blue nodes are not covered
- (b) Number of red nodes covered is minimized

- Greedy Algorithm:

- Construct “good” conjunctions of blocking predicates $\{p_1, p_2, \dots\}$.
- Pick k conjunctions $\{p_{i1}, p_{i2}, \dots, p_{ik}\}$, such that the following is minimized

$$\frac{\text{number of new blue nodes covered by } p_{ij}}{\text{number of red nodes covered by } p_{ij}}$$

Blocking Techniques

- Predicate blocking
- **Locality Sensitive Hashing**
- Canopy Clustering

minHash (Minwise Independent Permutations)

- Let F_x be a set of features for mention x
 - (predicates of) attribute values
 - character ngrams
 - optimal blocking functions ...
- Let π be a random permutation of features in F_x
 - E.g., order imposed by a random hash function
- $\text{minHash}(x)$ = minimum element in F_x according to π

Why minHash?

Surprising property: For a random permutation π ,

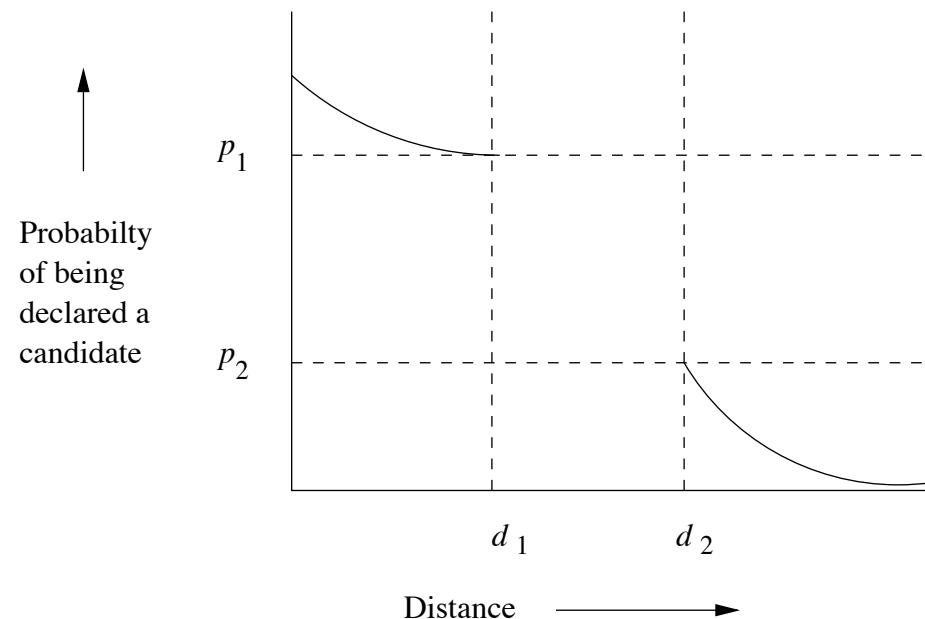
$$P[minHash(x) = minHash(y)] = \frac{F_x \cap F_y}{F_x \cup F_y}$$

Locality Sensitive Hashing Functions

Suppose d is a distance metric on a domain.

A family of functions \mathbf{F} is said to be (d_1, d_2, p_1, p_2) -sensitive if for all f in \mathbf{F} ,

- If $d(x, y) < d_1$,
then $P[f(x) = f(y)] > p_1$
- If $d(x, y) > d_2$,
then $P[f(x) = f(y)] < p_2$



Locality sensitive family for Jaccard

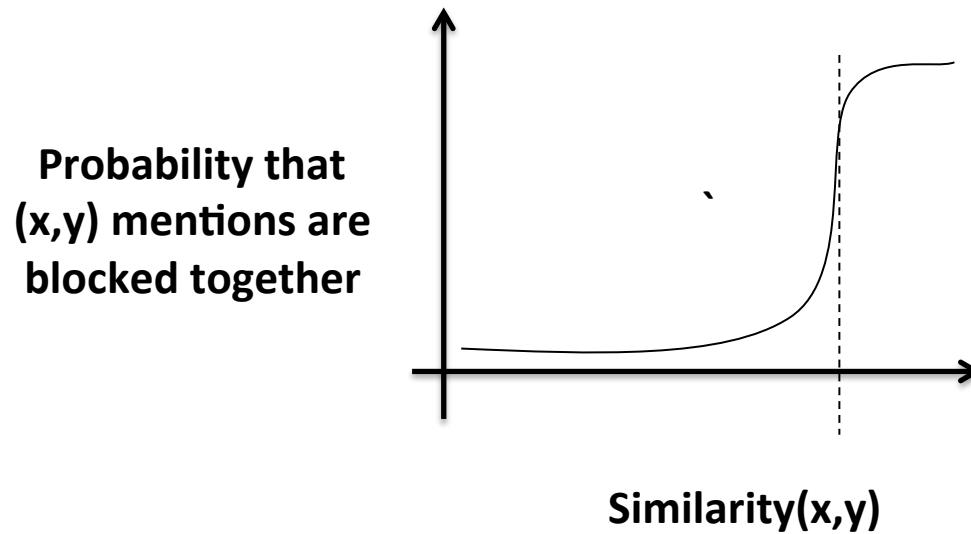
- Jaccard distance = 1 - Jaccard similarity = $1 - \frac{F_x \cap F_y}{F_x \cup F_y}$
- minHash is one example of locality sensitive family that can strongly distinguish pairs that are close from pairs that are far.
- The family of minHash functions is a $(d_1, d_2, 1-d_1, 1-d_2)$ -sensitive family for the Jaccard distance.

Blocking based on minHash

Surprising property: For a random permutation π ,

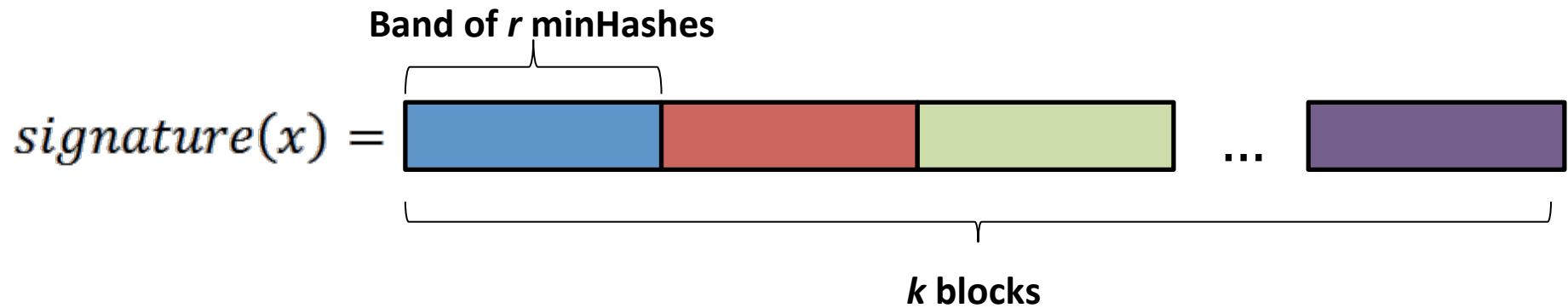
$$P[minHash(x) = minHash(y)] = \frac{F_x \cap F_y}{F_x \cup F_y}$$

How to build a blocking scheme such that only pairs with
Jacquard similarity $> s$ fall in the same block (with high prob)?



Blocking using minHashes

- Compute minHashes using $r * k$ permutations (hash functions)



- Signature's that match on ***1 out of k*** bands, go to the same block.

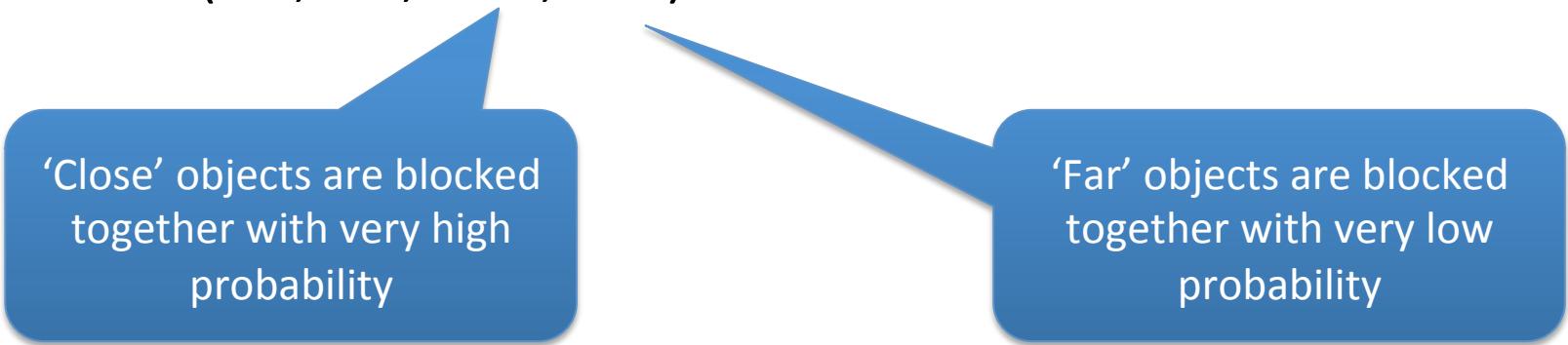
minHash Analysis

- Let F be a $(0.2, 0.6, 0.8, 0.4)$ -sensitive family of minHash functions
 - Pairs with Jaccard similarity > 0.8 are close, and similarity < 0.4 are far
- Let F_1 be the family constructed using a “band of $r=5$ minHashes” (**AND construction on F**)
 - F_1 is $(0.2, 0.6, 0.8^5, 0.4^5) = (0.2, 0.6, 0.328, 0.01)$ -sensitive

‘Far’ objects are blocked together with very low probability

minHash Analysis

- F is $(0.2, 0.6, 0.8, 0.4)$ -sensitive minHash
- F_1 is a “band of $r=5$ minHashes” (**AND construction on F**)
 - F_1 is $(0.2, 0.6, 0.8^5, 0.4^5) = (0.2, 0.6, 0.328, 0.01)$ -sensitive
- Let F_2 be the family constructed using “ $k=20$ bands of $r=5$ minHashes each” (**OR construction on F_1**)
 - F_2 is $(0.2, 0.6, 1 - (1 - 0.8^5)^{10}, 1 - (1 - 0.4^5)^{10})$
 $= (0.2, 0.6, 0.98, 0.09)$ -sensitive



‘Close’ objects are blocked together with very high probability

‘Far’ objects are blocked together with very low probability

minHash Analysis

- F is minHash family
 - $(0.2, 0.6, 0.8, 0.4)$ -sensitive
- F_1 is a “band of $r=5$ minHashes”
 - $(0.2, 0.6, 0.328, 0.01)$ -sensitive
- F_2 is “ $k = 20$ bands of $r=5$ minHashes each”
 - $(0.2, 0.6, 0.98, 0.09)$ -sensitive

$$r = 5, k = 20$$

Sim(s)	P(not same block)
0.9	0.9986
0.8	0.98
0.7	0.84
0.6	0.55
0.5	0.27
0.4	0.09
0.3	0.02
0.2	0.003
0.1	0.00009

LSH for Hamming distance

- Given two vectors x, y
- Hamming distance $h(x,y) = \text{number of positions where } x \text{ and } y \text{ are different}$
- minHash: $(d_1, d_2, 1-d_1/d, 1-d_2/d)$ -sensitive

LSH for Cosine Distance

- Cosine Distance: angle between two vectors
- Locality sensitive function F :

Generate v in $\{-1, +1\}^d$ (d is the dimensionality of x)
 $f(x) = f(y)$ if $x.vf$ and $y.vf$ have the same sign.

- F is $(d_1, d_2, (180-d_1)/180, d_2/180)$ -sensitive

Summary of Hash-based Blocking

- Complex boolean functions can be built to optimize recall using a training set of matches and non-matches
- Locality sensitive hashing functions can strongly distinguish pairs that are close from pairs that are far.
- AND and OR construction help amplify the distinguishing capability of locality sensitive functions.
- Can design good LSH family for many distance metrics.

Blocking Techniques

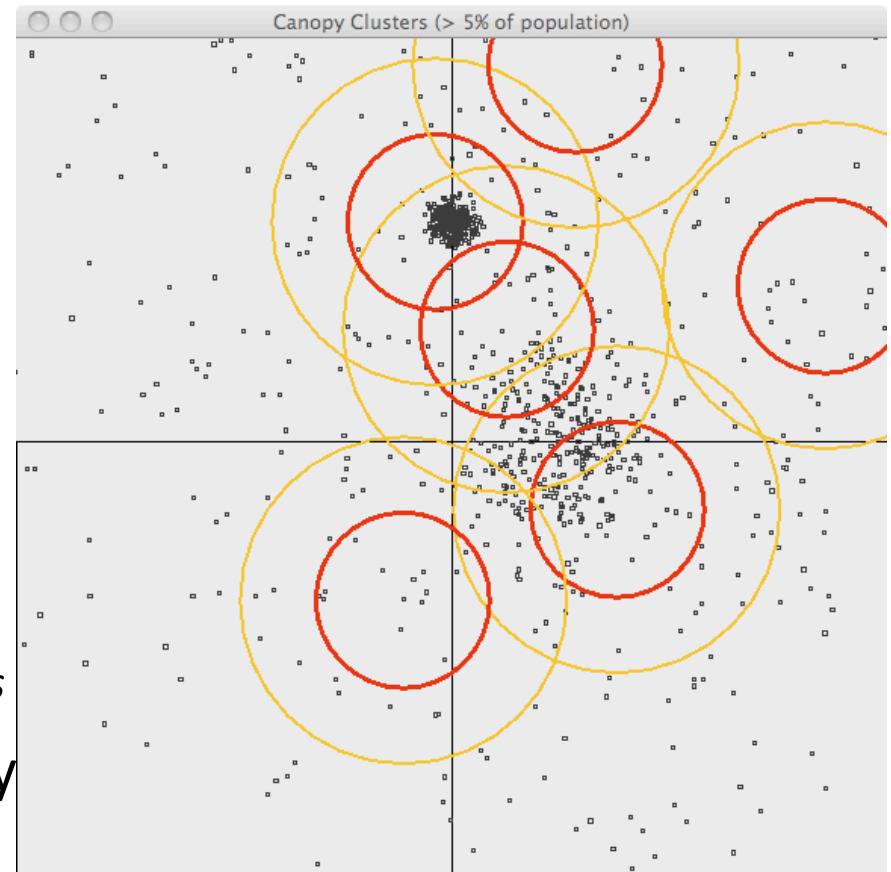
- Predicate blocking
- Locality Sensitive Hashing
- **Canopy Clustering**

Canopy Clustering [McCallum et al KDD'00]

Input: Mentions M ,
 $d(x,y)$, a distance metric,
thresholds $T_1 > T_2$

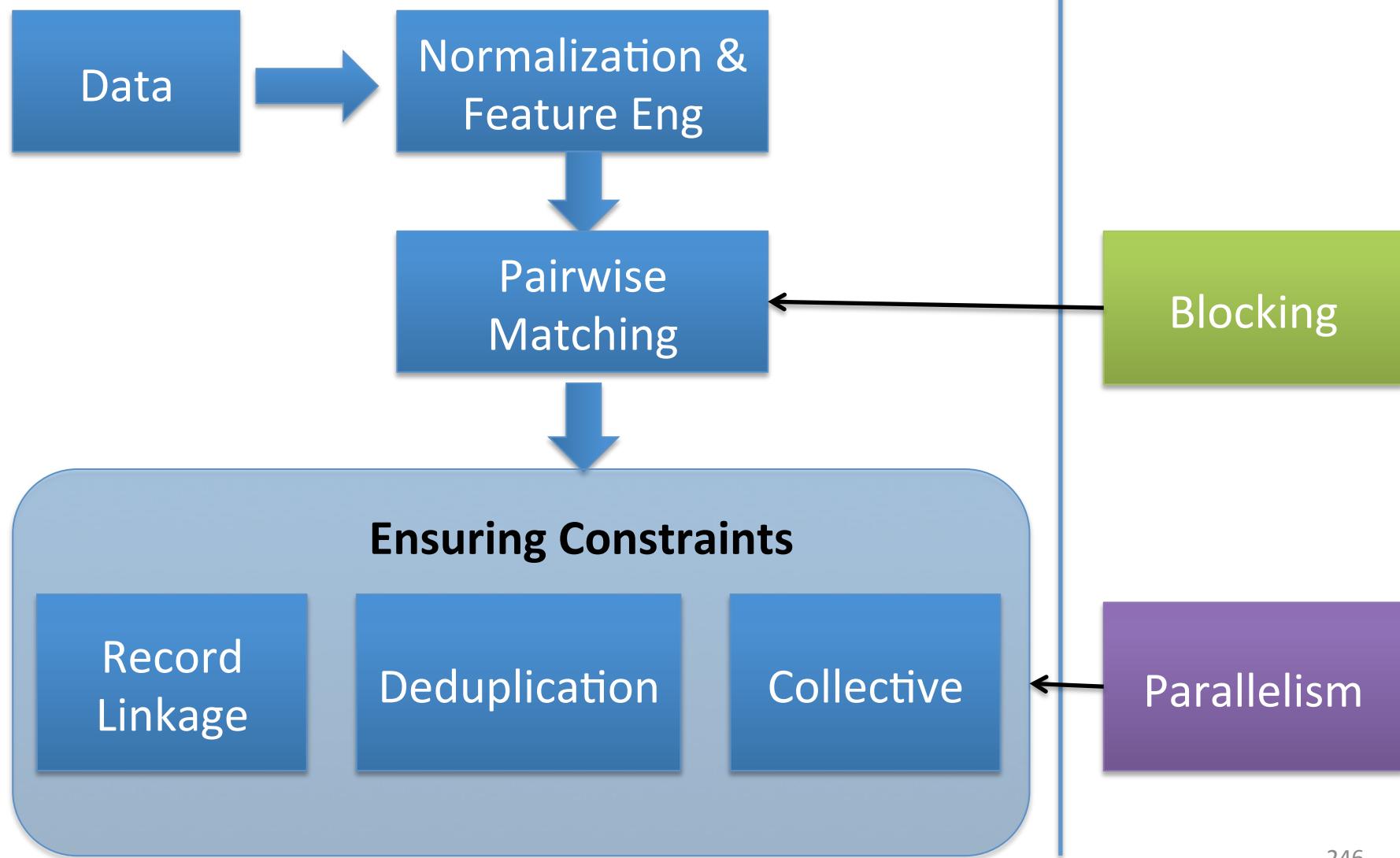
Algorithm:

1. Pick a random element x from M
2. Create new canopy C_x using
mentions y s.t. $d(x,y) < T_1$
3. Delete all mentions y from M
s.t. $d(x,y) < T_2$ (*from consideration in this*
4. Return to Step 1 if M is not empty



Part 2: Algorithmic Foundations

Part 3: Scaling

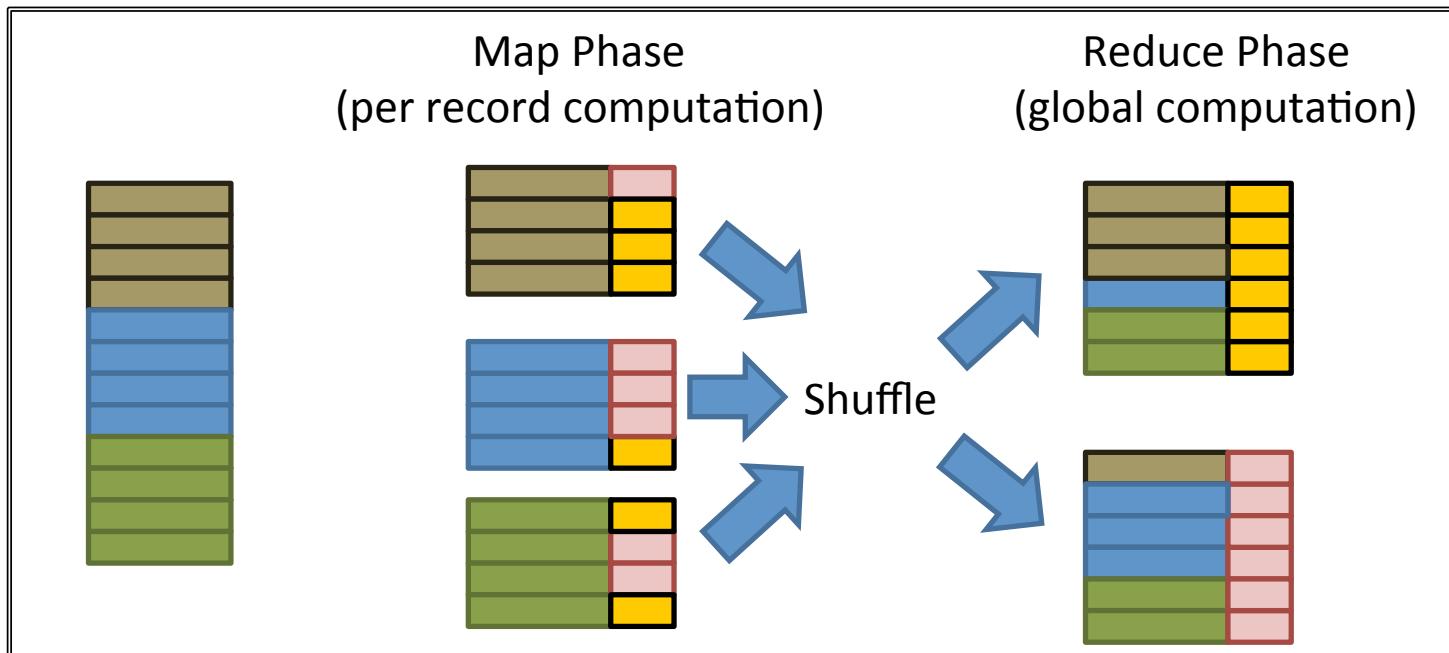


Distributed ER

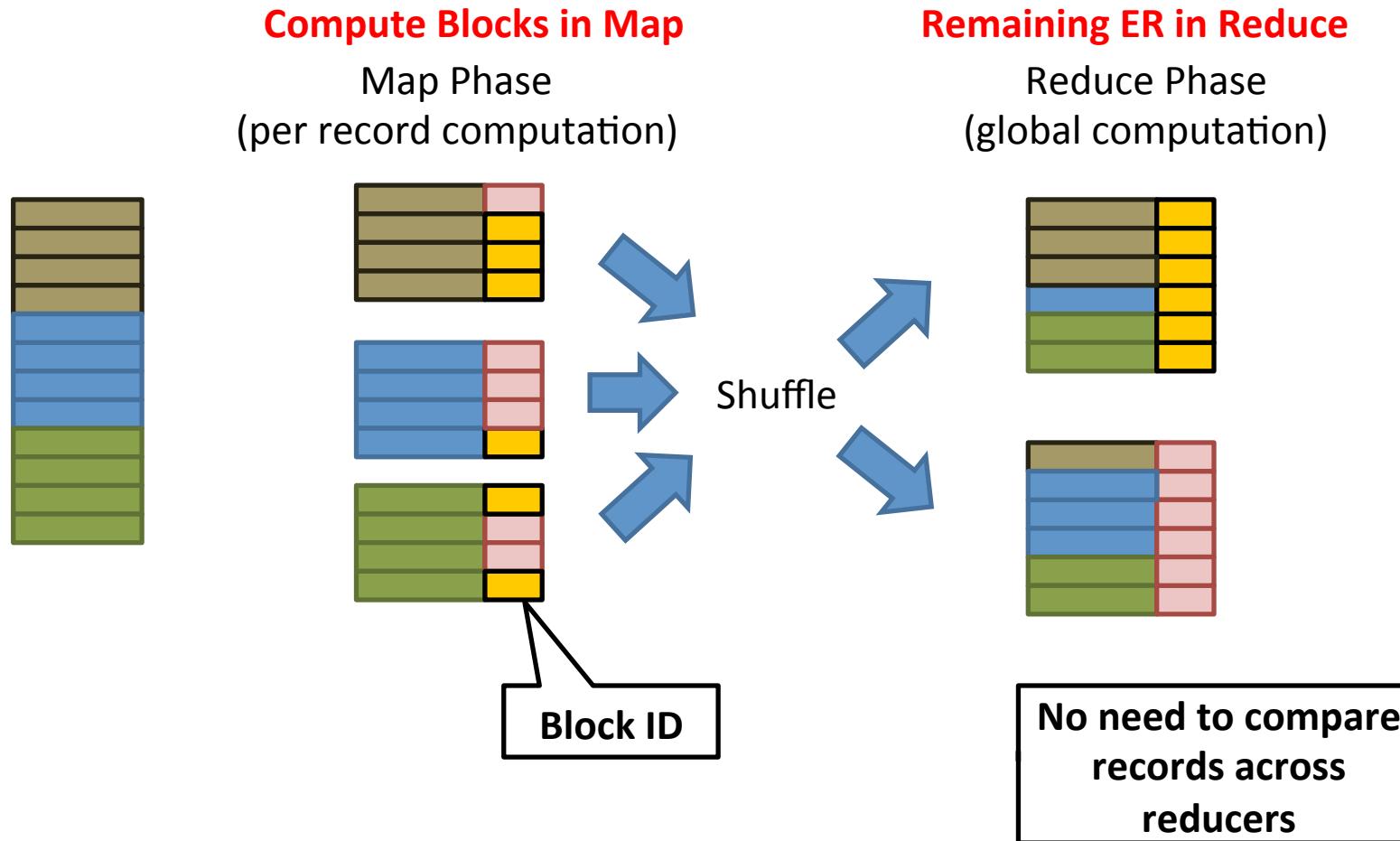
- Map-reduce is very popular for large tasks
 - Simple programming model for massively distributed data

```
map   (k1,v1) → list(k2,v2);  
reduce (k2,list(v2)) → list(k3,v3).
```

- Hadoop provides fault tolerance and is open source



ER with Disjoint Blocking



Non-disjoint Blocking

- How to block?
 - Hash-based: need an efficient technique to group records if they match on n -out-of- k blocking keys [Vernica et al SIGMOD'10]
 - Distance-based: canopy clustering on map-reduce [Mahout]
 - Iterative Blocking [Whang et al SIGMOD '09]

Problem: Information needed for a record is in multiple reducers.

- Information needed for a record is in multiple reducers.
 - Example 1:
 - Reducer 1: “a” matches with “b”
 - Reducer 2: “a” matches with “c”
 - Need to communicate in order to correctly resolve “a”, “b”, “c”

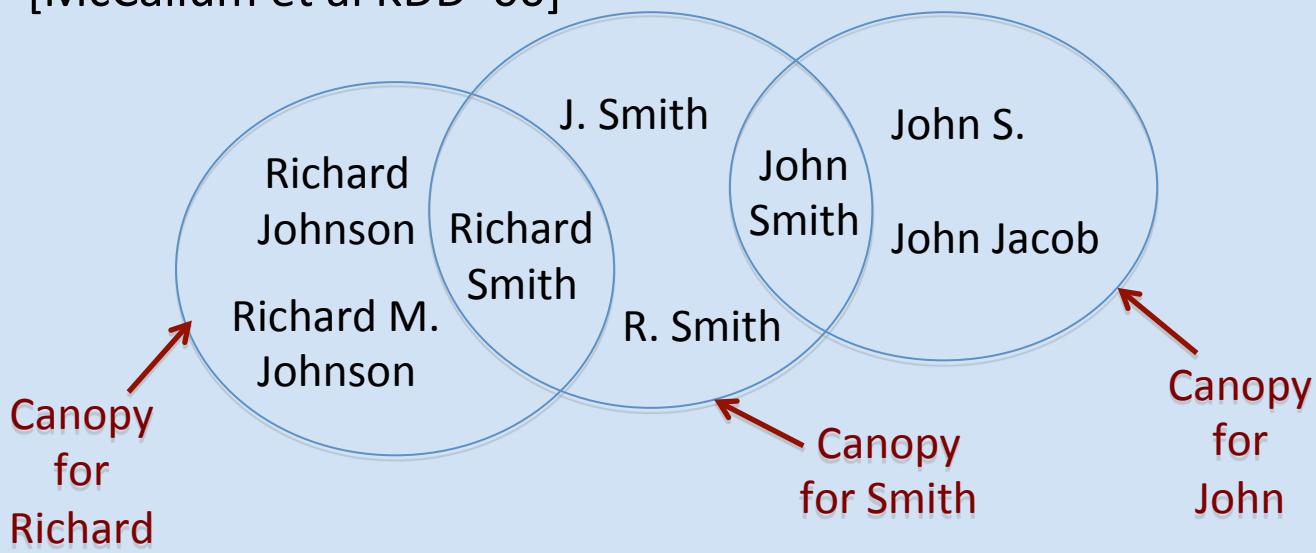
Problem: Information needed for a record is in multiple reducers.

Example 2: Dedup papers and authors

Id	Author-1	Author-2	Paper
A ₁	John Smith	Richard Johnson	Indices and Views
A ₂	J Smith	R Johnson	SQL Queries
A ₃	Dr. Smyth	R Johnson	Indices and Views

Problem: Information needed for a record is in multiple reducers.

Canopy clustering results in non-disjoint clusters
[McCallum et al KDD '00]



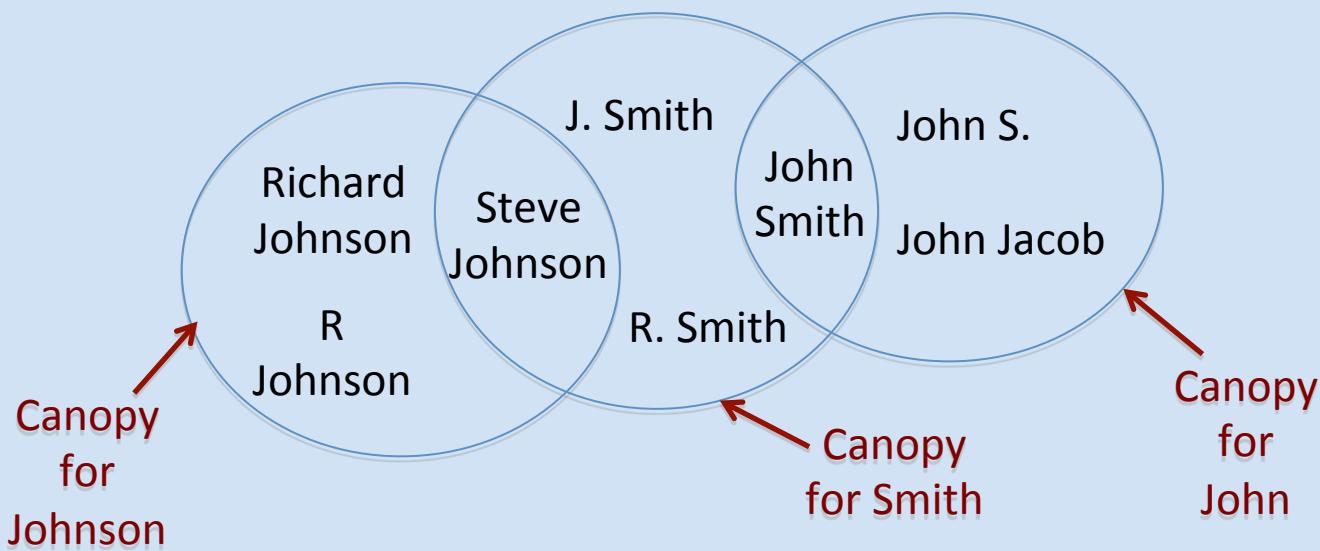
Slide adapted from [Rastogi et al VLDB11] talk

Problem: Information needed for a record is in multiple reducers.

$\text{CoAuthor}(A_1, B_1) \wedge \text{CoAuthor}(A_2, B_2) \wedge \text{match}(B_1, B_2) \rightarrow \text{match}(A_1, A_2)$

CoAuthor rule grounds to the correlation

$\text{match}(\text{Richard Johnson}, \text{R Johnson}) \Rightarrow \text{match}(\text{J. Smith}, \text{John Smith})$



253

Slide adapted from [Rastogi et al VLDB11] talk

Problem: Information needed for a record is in multiple reducers.

Solution 1: Efficiently find Connected Components [Rastogi et al 2013,
Kang et al ICDM 2009]

+ Correlation Clustering / Collective ER in each component

Solution 2: Correlation Clustering / Collective ER in each canopy

+ Message Passing [Rastogi et al VLDB'11]

Problem: Information needed for a record is in multiple reducers.

Solution 1: Efficiently find Connected Components [Rastogi et al 2012,
Kang et al ICDM 2009]
+ Correlation Clustering / Collective ER in each component

Connected components can be large in relational/multi-entity ER.

Solution 2: Correlation Clustering / Collective ER in each canopy
+ Message Passing [Rastogi et al VLDB'11]

Message Passing

Simple Message Passing (SMP)

1. Run entity matcher M locally in each canopy
2. If M finds a $\text{match}(r_1, r_2)$ in some canopy, pass it as evidence to all canopies
3. Rerun M within each canopy using new evidence
4. Repeat until no new matches found in each canopy

Runtime: $O(k^2 f(k) c)$

- k : maximum size of a canopy
- $f(k)$: Time taken by ER on canopy of size k
- c : number of canopies

Formal Properties

for a well behaved ER method ...

Convergence: No. of steps \leq no. of matches

Consistency: Output independent of the canopy order

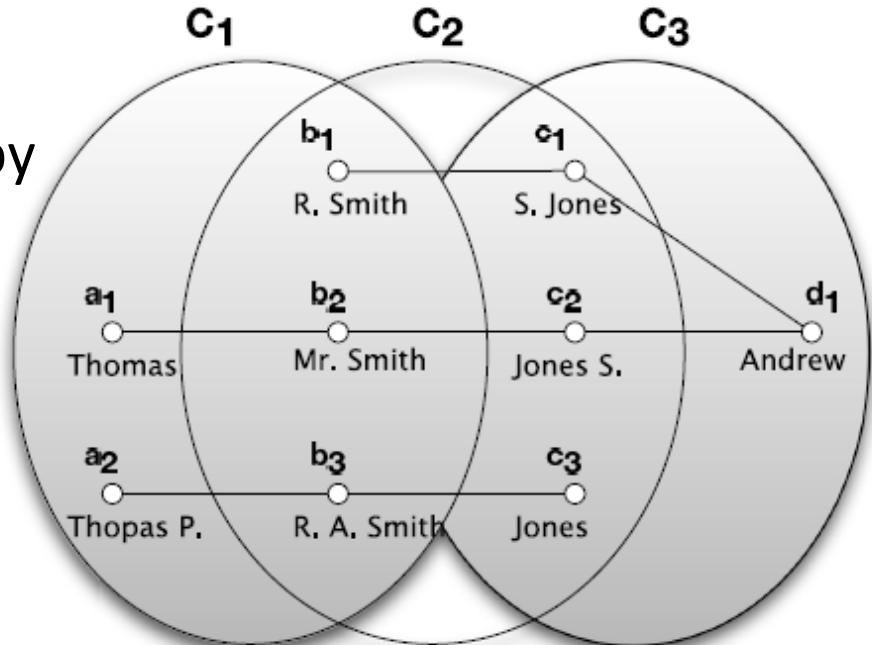
Soundness: Each output match is actually a true match

Completeness. Each true match is also a output match

Completeness

Papers 2 and 3 match only if a canopy knows that

- $\text{match}(a_1, a_2)$
- $\text{match}(b_2, b_3)$
- $\text{match}(c_2, c_3)$



Simple message passing will not find any matches
- thus, no messages are passed, no progress

Solution: Maximal message passing

- Send a message if there is a potential for match

Summary of Scalability

- $O(|R|^2)$ pairwise computations can be prohibitive.
 - Blocking eliminates comparisons on a large fraction of non-matches.
- Equality-based Blocking:
 - Construct (one or more) blocking keys from features
 - Records not matching on any key are not compared.
- Neighborhood based Blocking:
 - Form overlapping canopies of records based on similarity.
 - Only compare records within a cluster.
- Computing connected components/Message Passing in addition to blocking can help distribute ER.

Part 4

CHALLENGES AND FUTURE DIRECTIONS

Challenges (1)

So far, we have viewed ER as a one-time process applied to entire database; none of these hold in real world.

- Temporal ER
 - ER algorithms need to account for change in real world
 - Reasoning about multiple sources [Pal et al. WWW 12]
 - Model transitions [Li et al VLDB11]
- Reasoning about source quality
 - Sources are not independent
 - Copying Problem [Dong et al VLDB09]
- Query Time ER
 - How do we selectively determine the smallest number of records to resolve, so we get accurate results for a particular query?
 - Collective resolution for queries [Bhattacharya & Getoor JAIR07]

Challenges (2)

User interact with ER systems in a variety of ways

- Active Learning & Crowdsourcing
 - State of the art active learning techniques optimize 0-1 loss (false positive + false negative rate)
 - Wrong optimization metric for ER [Bellare et al KDD 2012]
 - Need to account for noisy training data from crowdsourcing
- ER & User-generated data
 - Deduplicated entities interact with users in the real world
 - Users tag/associate photos/reviews with businesses on Google / Yahoo
 - What should be done to support interactions?

Challenges (2)

User interact with ER systems in a variety of ways

- Active Learning
- ER & User-generated data
- UI to support ER
 - How to allow users to make ER decisions and see the context of their decisions?
 - D-Dupe, tool for relational ER
 - [Kang et al. TVCG 2008]
 - <http://www.cs.umd.edu/projects/linqs/ddupe>

Challenges (3)

Identity is not always a ‘crisp’ concept

- E.g. products: movies and books versus apparel and food
- E.g. composite objects: events, organizations, etc.

Open Issues

- ER is often part of bigger inference problem
 - Pipelined approaches and joint approaches to information extraction and graph identification
 - How can we characterize how ER errors affect overall quality of results?
- ER Theory
 - Need better support for theory which can give relational learning bounds
- ER & Privacy
 - ER enables record re-identification
 - How do we develop a theory of privacy-preserving ER?
- ER Benchmarks
 - Need for large-scale real-world ER datasets with groundtruth
 - Synthetic data useful for scaling but hard to capture rich complexities of real world

Summary

- We've covered algorithmic foundations for ER and touched on some of the open problems
- Growing omnipresence of massive linked data, and the need for creating knowledge bases from text and unstructured data motivate a number of new big data challenges in ER
- As data, noise, and knowledge grows, greater needs & opportunities for intelligent reasoning about entity resolution

THANK YOU!

References – Intro

- W. Willinger et al, “Mathematics and the Internet: A Source of Enormous Confusion and Great Potential”, Notices of the AMS 56(5), 2009
- L. Gill and M. Goldcare, “English National Record Linkage of Hospital Episode Statistics and Death Registration Records”, Report to the Department of Health, 2003
- T. Herzog et al, “Data Quality and Record Linkage Techniques”, Springer 2007
- A. Elmagrid et al, “Duplicate Record Detection”, TKDE 2007
- P. Christen, “Data Matching”, Springer 2012
- N. Koudas et al, “Record Linkage: Similarity measures and Algorithms”, SIGMOD 2006
- X. Dong & F. Naumann, “Data fusion--Resolving data conflicts for integration”, VLDB 2009
- L. Getoor & A. Machanavajjhala, “Entity Resolution: Theory, Practice and Open Challenges”, AAAI 2012

References – Single Entity ER

- D. Menestrina et al, “Evaluation Entity Resolution Results”, PVLDB 3(1-2), 2010
- M. Cochiniwala et al, “Efficient data reconciliation”, Information Sciences 137(1-4), 2001
- M. Bilenko & R. Mooney, “Adaptive Duplicate Detection Using Learnable String Similarity Measures”, KDD 2003
- P. Christen, “Automatic record linkage using seeded nearest neighbour and support vector machine classification.”, KDD 2008
- Z. Chen et al, “Exploiting context analysis for combining multiple entity resolution systems”, SIGMOD 2009
- A. McCallum & B. Wellner, “Conditional Models of Identity Uncertainty with Application to Noun Coreference”, NIPS 2004
- H. Newcombe et al, “Automatic linkage of vital records”, Science 1959
- I. Fellegi & A. Sunter, “A Theory for Record Linkage”, JASA 1969
- W. Winkler, “Overview of Record Linkage and Current Research Directions”, Research Report Series, US Census, 2006
- T. Herzog et al, “Data Quality and Record Linkage Techniques”, Springer, 2007
- P. Ravikumar & W. Cohen, “A Hierarchical Graphical Model for Record Linkage”, UAI 2004

References – Single Entity ER (contd.)

- S. Sarawagi et al, “Interactive Deduplication using Active Learning”, KDD 2000
- S. Tejada et al, “Learning Object Identification Rules for Information Integration”, IS 2001
- A. Arasu et al, “On active learning of record matching packages”, SIGMOD 2010
- K. Bellare et al, “Active sampling for entity matching”, KDD 2012
- A. Beygelzimer et al, “Agnostic Active Learning without Constraints”, NIPS 2010
- J. Wang et al, “CrowdER: Crowdsourcing Entity Resolution”, PVLDB 5(11), 2012
- A. Marcus et al, “Human-powered Sorts and Joins”, PVLDB 5(1), 2011

References – Single Entity ER (contd.)

- R. Gupta & S. Sarawagi, “Answering Table Augmentation Queries from Unstructured Lists on the Web”, PVLDB 2(1), 2009
- A. Das Sarma et al, “An Automatic Blocking Mechanism for Large-Scale De-duplication Tasks”, CIKM 2012
- M. Bilenko et al, “Adaptive Product Normalization: Using Online Learning for Record Linkage in Comparison Shopping”, ICDM 2005
- S. Chaudhuri et al, “Robust Identification of Fuzzy Duplicates”, ICDE 2005
- W. Soon et al, “A machine learning approach to coreference resolution of noun phrases”, Computational Linguistics 27(4) 2001
- N. Bansal et al, “Correlation Clustering”, Machine Learning 56(1-3), 2004
- V. Ng & C. Cardie, “Improving machine learning approaches to coreference resolution”, ACL 2002
- M. Elsner & E. Charnaiik, “You talking to me? a corpus and algorithm for conversation disentanglement”, ACL-HLT 2008
- M. Elsner & W. Schudy, “Bounding and Comparing Methods for Correlation Clustering Beyond ILP”, ILP-NLP 2009
- N. Ailon et al, “Aggregating inconsistent information: Ranking and clustering”, JACM 55(5), 2008
- X. Dong et al, “Integrating Conflicting Data: The Role of Source Dependence”, PVLDB 2(1), 2009
- A. Pal et al, “Information Integration over Time in Unreliable and Uncertain Environments”, WWW 2012
- A. Culotta et al, “Canonicalization of Database Records using Adaptive Similarity Measures”, KDD 2007
- O. Benjelloun et al, “Swoosh: A generic approach to Entity Resolution”, VLDBJ 18(1), 2009

References – Constraints & Multi-Relational ER

- R. Ananthakrishna et. al, "Eliminating fuzzy duplicates in data warehouses", VLDB 2002
- A. Arasu et al, "Large-Scale Deduplication with Constraints using Dedupalog", ICDE 2009
- S. Chaudhuri et al., "Leveraging aggregate constraints for deduplication", SIGMOD07
- X. Dong et al, "Reference Recconciliation in Complex Information Spaces", SIGMOD 2005
- I. Bhattacharya & L. Getoor, "Collective Entity Resolution in Relational Data", TKDD 2007
- I. Bhattacharya & L. Getoor, "A Latent Dirichlet Model for Unsupervised Entity Resolution ", SDM 2007
- S. Sarawagi & W. Cohen, "Semi-Markov conditional random fields for information extraction", NIPS, 2004.
- P. Bohannon et al., "Conditional Functional Dependencies for Data Cleaning", ICDE 2007
- M. Broeckeler & L. Getoor , "Probabilistic Similarity Logic", UAI 2010
- W. Fan, "Dependencies revisited for improving data quality", PODS 2008
- H. Pasula et al , "Identity Uncertainty and Citation Matching", NIPS 2002
- D. Kalashnikov et al, "Domain-Independent Data Cleaning via Analysis of Entity-Relationship Graph", TODS06
- J. Lafferty et al, "Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data.", ICML 2001
- X. Li et al, "Identification and Tracing of Ambiguous Names: Discriminative and Generative Approaches", AAAI 2004
- A. McCallum & B. Wellner, "Conditional Models of Identity Uncertainty with Application to Noun Coreference", NIPS 2004
- M. Richardson & P. Domingos, "Markov Logic", Machine Learning 62, 2006
- W. Shen et al., "Constraint-based Entity Matching", AAAI 2005
- P. Singla & P. Domingos, "Entity Resolution with Markov Logic", ICDM 2006
- M. Wick, A. Culotta, K. Rohanimanesh, A. McCallum, " An Entity Based Model for Coreference Resolution", SDM 2009
- Whang et al., "Generic Entity Resolution with Negative Rules", VLDBJ 2009
- Whang et al., "Joint Entity Resolution", ICDE 2012

References – Blocking

- M. Bilenko et al, “Adaptive Blocking: Learning to Scale Up Record Linkage and Clustering”, ICDM 2006
- M. Michelson & C. Knoblock, “Learning Blocking Schemes for Record Linkage”, AAAI 2006
- A. Das Sarma et al, “An Automatic Blocking Mechanism for Large-Scale De-duplication Tasks”, CIKM 2012
- A. Broder et al, “Min-Wise Independent Permutations”, STOC 1998
- G. Papadias et al, “Beyond 100 million entities: large-scale blocking-based resolution for heterogenous data,” WSDM 2012
- M. Hernandez & S. Stolfo, “The merge/purge problem for large databases”, SIGMOD 1995
- A. McCallum et al, “Efficient clustering of high-dimensional data sets with application to reference matching”, KDD 2000
- L. Kolb et al, “Dedoop: Efficient deduplication with Hadoop”, (demo) PVLDB 5(12), 2012
- R. Vernica et al, “Efficient Parallel Set-Similarity Joins Using MapReduce”, SIGMOD 2010
- Apache Mahout: Scalable Machine Learning and Data Mining, <http://mahout.apache.org/>
- S. Whang et al, “Entity Resolution with Iterative Blocking”, SIGMOD 2009
- U. Kang et al, “PEGASUS: A Peta-Scale Graph Mining System - Implementation and Observations”, ICDM 2009
- V. Rastogi et al, “Finding Connected Components on Map-reduce in Poly-Log Rounds”, Corr 2012
- V. Rastogi et al, “Large-Scale Collective Entity Matching”, PVLDB 4(4), 2011

References – Challenges & Future Directions

- I. Bhattacharya and L. Getoor, "Query-time Entity Resolution", JAIR 2007
- X. Dong, L. Berti-Equille, D. Srivastava, "Truth discovery and copying detection in a dynamic world", VLDB 2009
- P. Li, X. Dong, A. Maurino, D. Srivastava, "Linking Temporal Records", VLDB 2011
- A. Pal, V. Rastogi, A. Machanavajjhala, P. Bohannon, "Information integration over time in unreliable and uncertain environments", WWW 2012