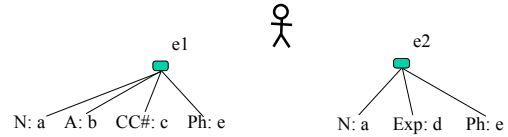


Generic Entity Resolution: Identifying Real-World Entities in Large Data Sets

Steven Whang
Stanford University

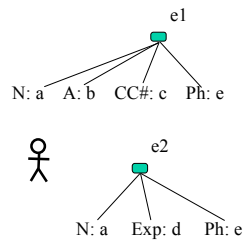
Work with: Hector Garcia-Molina, David Menestrina, Omar Benjelloun, Qi Su, Jennifer Widom, Tyson Condie, Nicolas Pombourcq, Tait Larson, Johnson Gong, Makoto Tachibana, Sutthipong Thavisomboon, Georgia Koutrika, Martin Theobald

Entity Resolution



Applications

- comparison shopping
- mailing lists
- classified ads
- customer files
- counter-terrorism

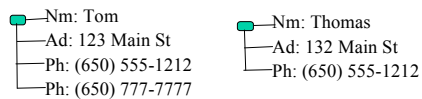


Outline

- Why is ER challenging?
- How is ER done?
- Some ER work at Stanford
 - Generic ER
 - Distributed ER
 - Iterative Blocking

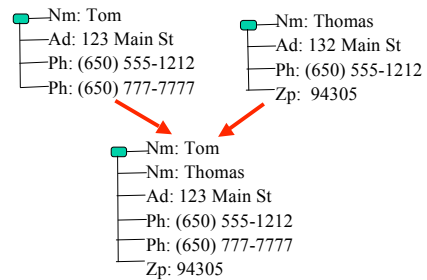
Challenges (1)

- No keys!
- Value matching
 - “Kaddafi”, “Qaddafi”, “Kadafi”, “Kaddaffi”...
- Record matching



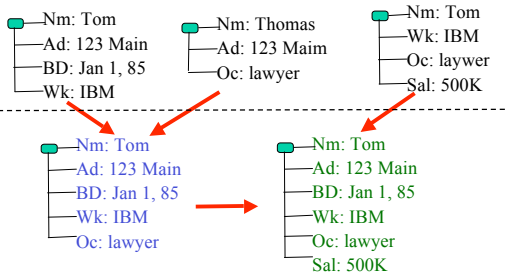
Challenges (2)

- Merging records



Challenges (3)

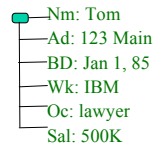
- Chaining



7

Challenges (4)

- Un-merging

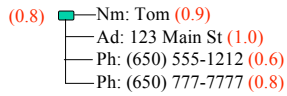


too young to make 500K at IBM!!

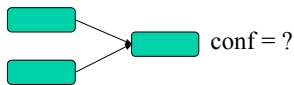
8

Challenges (5)

- Confidences in data



- In value matching, match rules, merge:



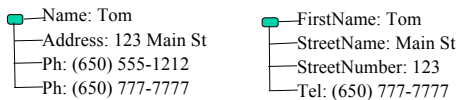
9

Taxonomy

- Pairwise snaps vs. clustering
- De-duplication vs. fidelity enhancement
- Schema differences
- Relationships
- Exact vs. approximate
- Generic vs application specific
- Confidences

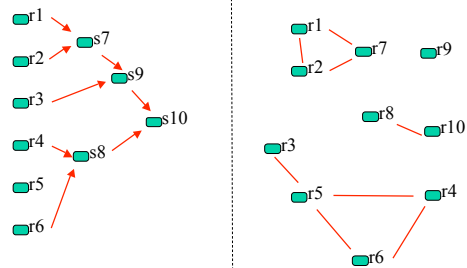
10

Schema Differences

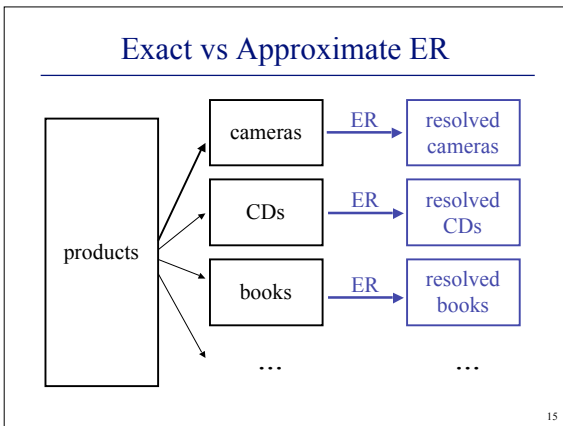
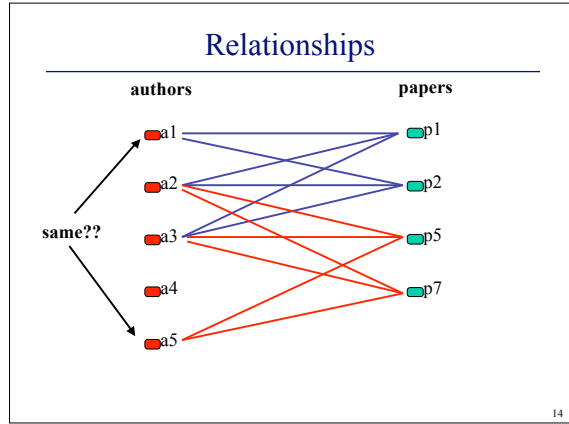
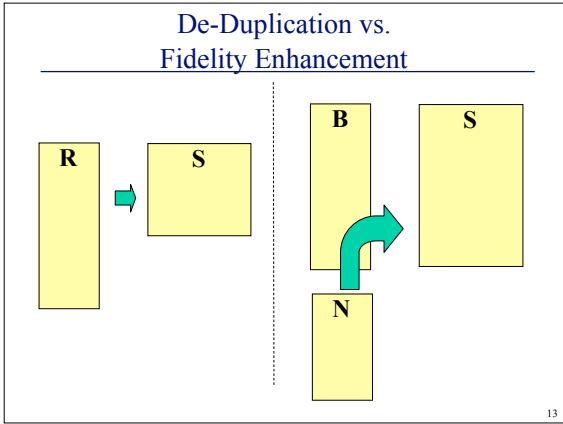


11

Pair-Wise Snaps vs. Clustering



12



- ### Generic vs Application Specific
- Match function $M(r, s)$
 - Merge function $\langle r, s \rangle \Rightarrow t$
- 16

- ### Taxonomy
- Pairwise snaps vs. clustering
 - De-duplication vs. fidelity enhancement
 - Schema differences
 - Relationships
 - Exact vs. approximate
 - Generic vs application specific
 - Confidences
- 17

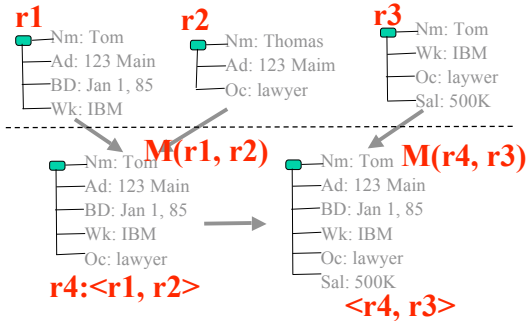
- ### Outline
- Why is ER challenging?
 - How is ER done?
 - Some ER work at Stanford
 - Generic ER ←
 - Distributed ER
 - Iterative Blocking
- 18

Taxonomy

- Pairwise snaps vs. clustering
- De-duplication vs. fidelity enhancement
- Schema differences No
- Relationships No
- Exact vs. approximate
- Generic vs application specific
- Confidences No

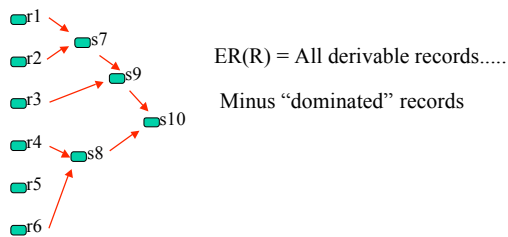
19

Model



20

Correct Answer



21

Question

- What is best sequence of match, merge calls that give us right answer?

22

Brute Force Algorithm

- Input R:
 - r1 = [a:1, b:2]
 - r2 = [a:1, c: 4, e:5]
 - r3 = [b:2, c:4, f:6]
 - r4 = [a:7, e:5, f:6]

23

Brute Force Algorithm

- Input R:
 - r1 = [a:1, b:2]
 - r2 = [a:1, c: 4, e:5]
 - r3 = [b:2, c:4, f:6]
 - r4 = [a:7, e:5, f:6]
- Match all pairs:
 - r1 = [a:1, b:2]
 - r2 = [a:1, c: 4, e:5]
 - r3 = [b:2, c:4, f:6]
 - r4 = [a:7, e:5, f:6]
 - r12 = [a:1, b:2, c:4, e:5]

24

Brute Force Algorithm

- Match all pairs:
 - r1 = [a:1, b:2]
 - r2 = [a:1, c: 4, e:5]
 - r3 = [b:2, c:4, f:6]
 - r4 = [a:7, e:5, f:6]
 - r12 = [a:1, b:2, c:4, e:5]
- Repeat:
 - r1 = [a:1, b:2]
 - r2 = [a:1, c: 4, e:5]
 - r3 = [b:2, c:4, f:6]
 - r4 = [a:7, e:5, f:6]
 - r12 = [a:1, b:2, c:4, e:5]
 - r123 = [a:1, b:2, c:4, e:5, f:6]

25

Question # 1

Brute Force Algorithm

- Match all pairs:
 - r1 = [a:1, b:2]
 - r2 = [a:1, c: 4, e:5]
 - r3 = [b:2, c:4, f:6]
 - r4 = [a:7, e:5, f:6]
 - r12 = [a:1, b:2, c:4, e:5]
- Repeat:
 - r1 = [a:1, b:2]
 - r2 = [a:1, c: 4, e:5]
 - r3 = [b:2, c:4, f:6]
 - r4 = [a:7, e:5, f:6]
 - r12 = [a:1, b:2, c:4, e:5]
 - r123 = [a:1, b:2, c:4, e:5, f:6]

Can we avoid comparisons?

26

Question # 2

Brute Force Algorithm

- Input R:
 - r1 = [a:1, b:2]
 - r2 = [a:1, c: 4, e:5]
 - r3 = [b:2, c:4, f:6]
 - r4 = [a:7, e:5, f:6]
- Match all pairs:
 - r1 = [a:1, b:2]
 - r2 = [a:1, c: 4, e:5]
 - r3 = [b:2, c:4, f:6]
 - r4 = [a:7, e:5, f:6]
 - r12 = [a:1, b:2, c:4, e:5]

Can we delete r1, r2?

27

ICAR Properties

- Idempotence:
 - $M(r1, r1) = \text{true}; \langle r1, r1 \rangle = r1$
- Commutativity:
 - $M(r1, r2) = M(r2, r1)$
 - $\langle r1, r2 \rangle = \langle r2, r1 \rangle$
- Associativity
 - $\langle r1, \langle r2, r3 \rangle \rangle = \langle \langle r1, r2 \rangle, r3 \rangle$

28

More Properties

- Representativity
 - If $\langle r1, r2 \rangle = r3$, then for any r4 such that $M(r1, r4)$ is true we also have $M(r3, r4) = \text{true}$.

29

ICAR Properties ➔ Efficiency

- Commutativity
- Idempotence
- Associativity
- Representativity

• Can discard records
• ER result independent of processing order

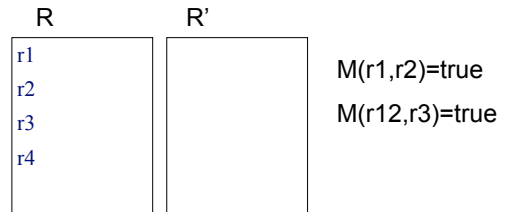
30

Swoosh Algorithms

- R-Swoosh
 - Merges records as soon as they match
 - Optimal in terms of record comparisons
- F-Swoosh
 - Remembers values seen for each feature
 - Avoids redundant value comparisons

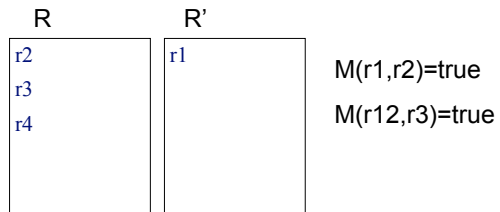
31

R-Swoosh



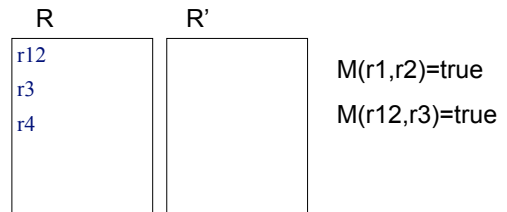
32

R-Swoosh



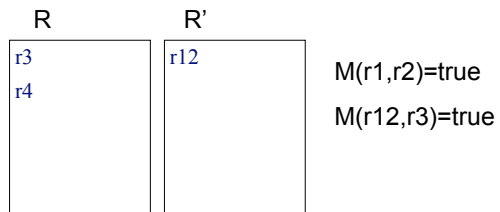
33

R-Swoosh



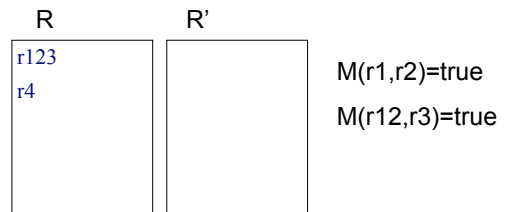
34

R-Swoosh

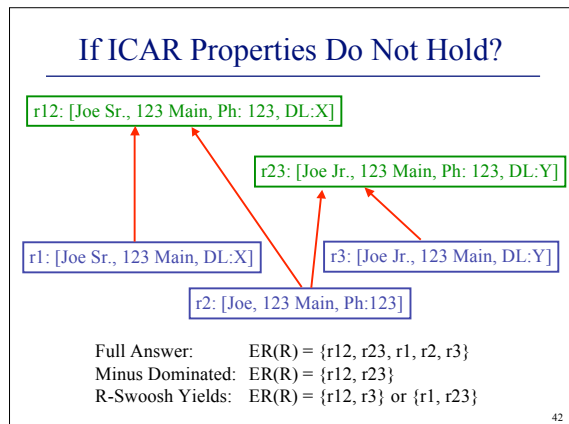
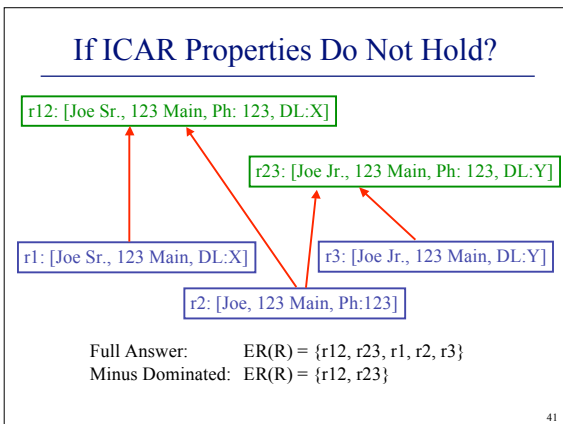
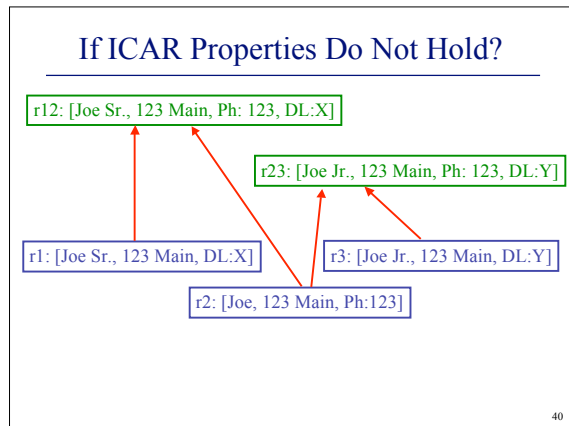
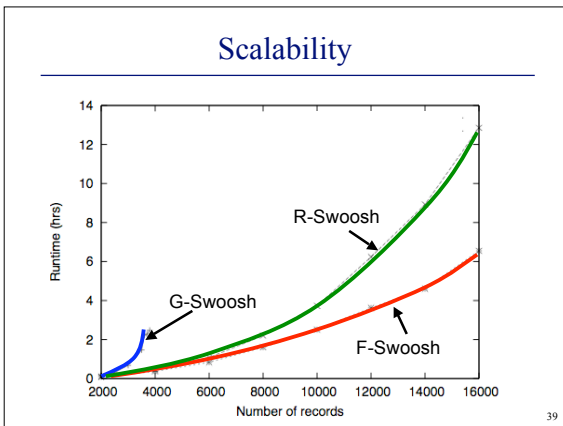
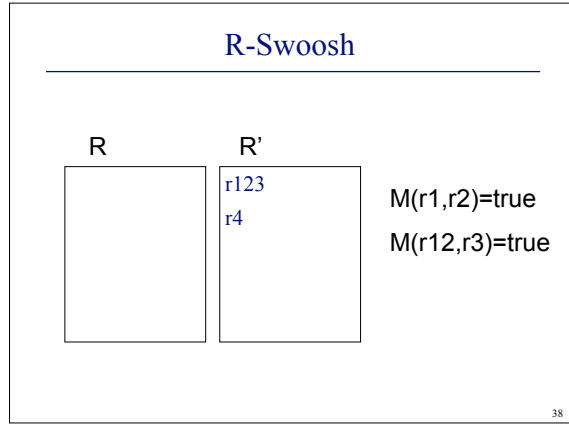
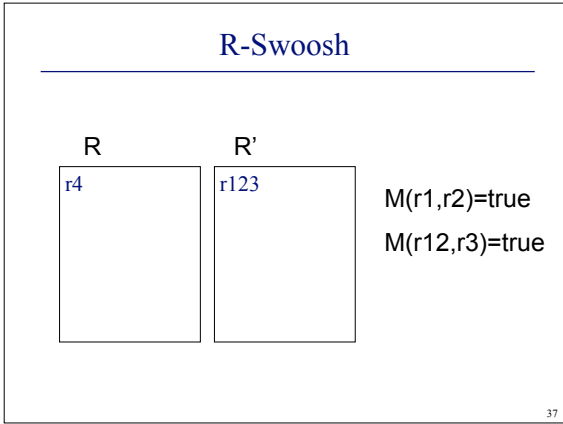


35

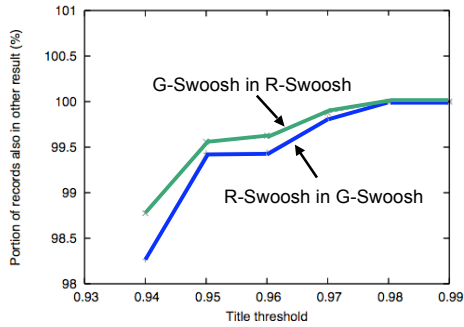
R-Swoosh



36



Swoosh Without ICAR Properties



43

Generic ER Summary

- Entity resolution is critical
- Generic approach yields reusable techniques
- Efficient resolution is important

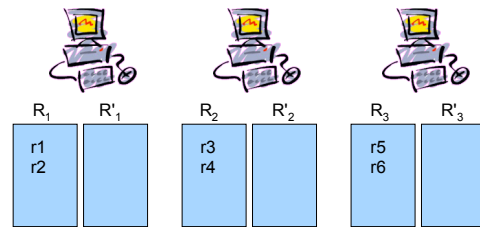
44

Outline

- Why is ER challenging?
- How is ER done?
- Some ER work at Stanford
 - Generic ER
 - Distributed ER
 - Iterative Blocking

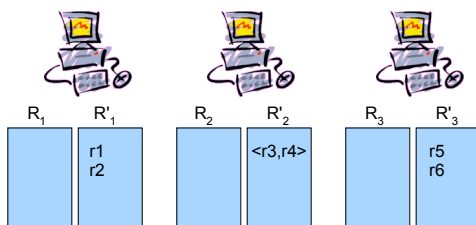
45

Partition Input Data



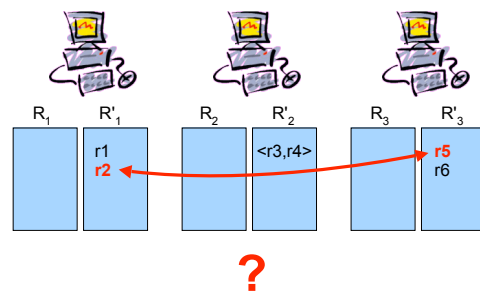
46

Partition Input Data



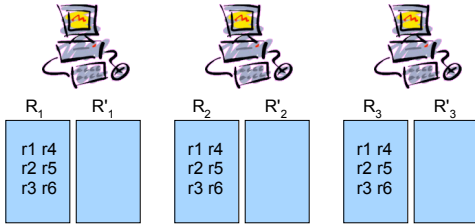
47

Partition Input Data



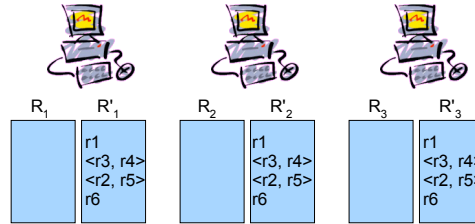
48

Partition Input Data



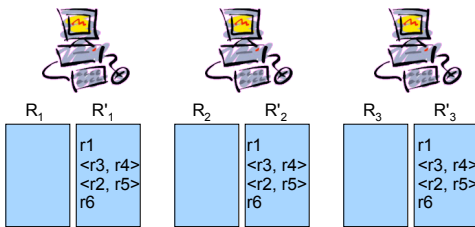
49

Partition Input Data



50

Partition Input Data



Redundant work!

51

Scope and Responsible

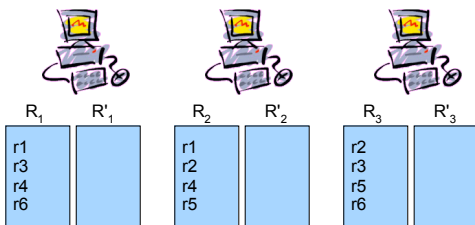
- $scope(r)$ – Returns list of processors that receive r
- $resp(p_k, r_i, r_j)$ – true if processor p_k should perform comparison of r_i and r_j

Coverage property:

For any pair of matching records r, r' , there exists at least one processor P_k such that $P_k \in scope(r) \cap scope(r')$ and $resp(P_k, r, r') = true$.

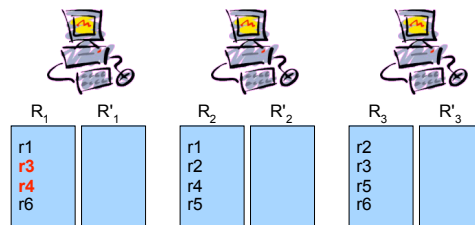
52

D-Swoosh

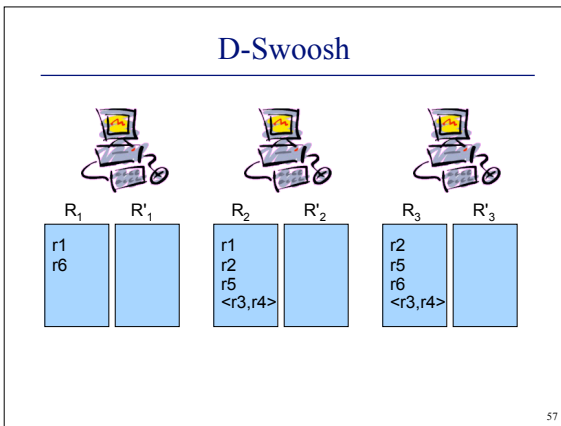
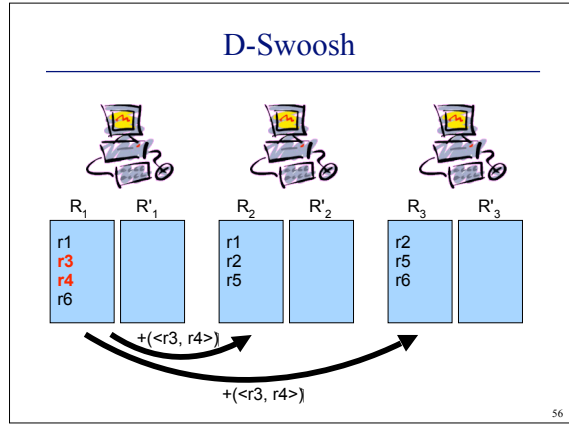
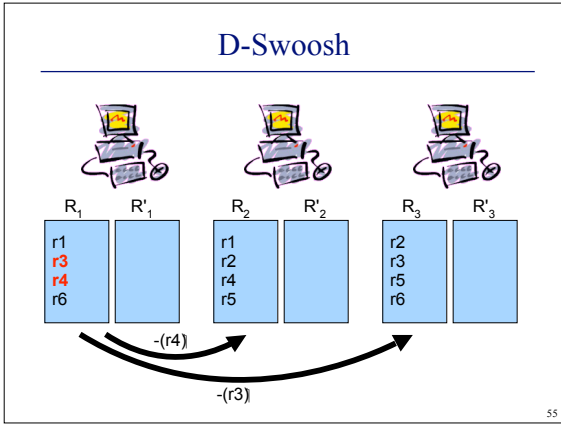


53

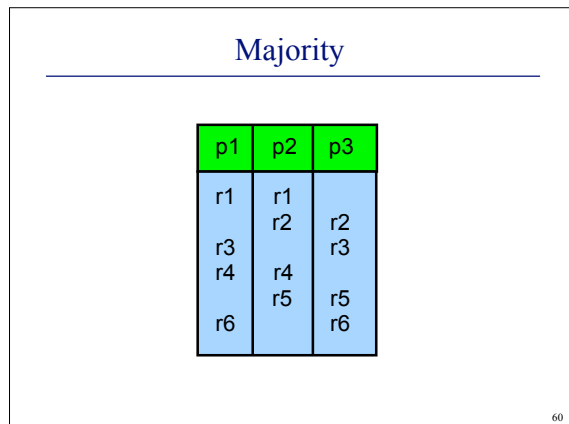
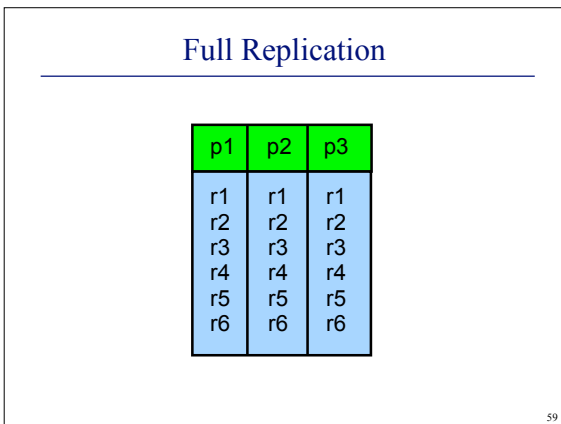
D-Swoosh

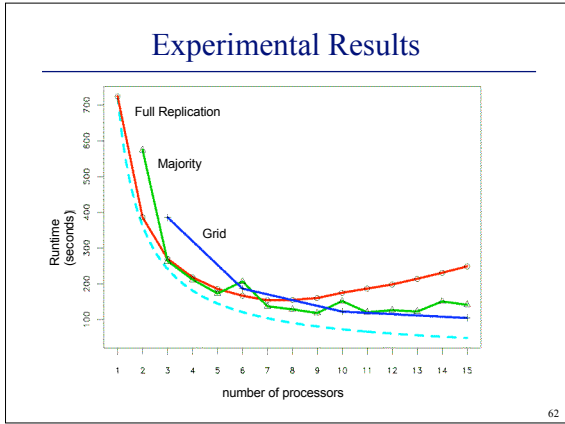
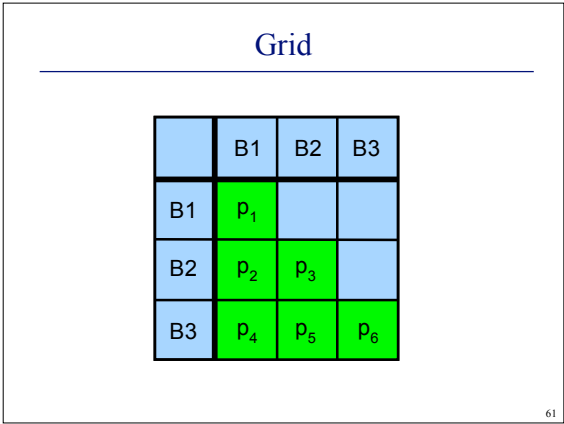


54



- ### Defining *scope* and *resp*
- What are good *scope* and *resp* functions?
 - How do different functions compare and scale?
 - How can we exploit semantic knowledge?
- 58





- ### Distributed ER Summary
- Entity Resolution is fundamentally expensive
 - Reduce processing time by:
 - Distributing data
 - Eliminating redundant work
 - ER benefits greatly from distributed computing
- 63

- ### Outline
- Why is ER challenging?
 - How is ER done?
 - Some ER work at Stanford
 - Generic ER
 - Distributed ER
 - Iterative Blocking ←
- 64

Entity Resolution

Record	Name	Zip	Email
r ₁	John Doe	12345	jdoe@yahoo
r ₂	John Doe	94305	
r ₃	J. Foe	94305	jdoe@yahoo

65

Entity Resolution

Record	Name	Zip	Email
r ₁	John Doe	12345	jdoe@yahoo
r ₂	John Doe	94305	
r ₃	J. Foe	94305	jdoe@yahoo
r ₁₂	John Doe	{12345, 94305}	jdoe@yahoo

r₁₂
↑ ↑
r₁ r₂ r₃

66

Entity Resolution

Record	Name	Zip	Email
r ₁	John Doe	12345	jdoe@yahoo
r ₂	John Doe	94305	
r ₃	J. Doe	94305	jdoe@yahoo
r ₁₂	John Doe	{12345, 94305}	jdoe@yahoo
r ₁₂₃	{John Doe, J. Doe}	{12345, 94305}	jdoe@yahoo

67

Entity Resolution

Record	Name	Zip	Email
r ₁	John Doe	12345	jdoe@yahoo
r ₂	John Doe	94305	
r ₃	J. Doe	94305	jdoe@yahoo
r ₁₂	John Doe	{12345, 94305}	jdoe@yahoo
r ₁₂₃	{John Doe, J. Doe}	{12345, 94305}	jdoe@yahoo

ER Solution: {r₁₂₃}

68

Blocking

Record	Name	Zip	Email
r ₁	John Doe	12345	jdoe@yahoo
r ₂	John Doe	94305	
r ₃	J. Doe	94305	jdoe@yahoo

Partition by	Block#1	Block#2
zip	r ₁	r ₂ , r ₃

69

Blocking

Record	Name	Zip	Email
r ₁	John Doe	12345	jdoe@yahoo
r ₂	John Doe	94305	
r ₃	J. Doe	94305	jdoe@yahoo

Partition by	Block#1	Block#2
zip	r ₁	r ₂ , r ₃
1 st char last name	r ₁ , r ₂	r ₃

70

Blocking

Record	Name	Zip	Email
r ₁	John Doe	12345	jdoe@yahoo
r ₂	John Doe	94305	
r ₃	J. Doe	94305	jdoe@yahoo

Partition by	Block#1	Block#2
zip	r ₁	r ₂ , r ₃
1 st char last name	r ₁₂	r ₃

71

Blocking

Record	Name	Zip	Email
r ₁	John Doe	12345	jdoe@yahoo
r ₂	John Doe	94305	
r ₃	J. Doe	94305	jdoe@yahoo

Partition by	Block#1	Block#2
zip	r ₁	r ₂ , r ₃
1 st char last name	r ₁₂	r ₃

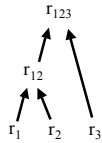
Blocking Solution: {r₁₂, r₃}

72

Iterative Blocking

Record	Name	Zip	Email
r ₁	John Doe	12345	jdoe@yahoo
r ₂	John Doe	94305	
r ₃	J. Doe	94305	jdoe@yahoo

Partition by	Block#1	Block#2
zip	r ₁	r ₂ , r ₃
1 st char last name	r ₁ , r ₂	r ₃

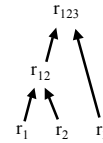


73

Iterative Blocking

Record	Name	Zip	Email
r ₁	John Doe	12345	jdoe@yahoo
r ₂	John Doe	94305	
r ₃	J. Doe	94305	jdoe@yahoo

Partition by	Block#1	Block#2
zip	r ₁	r ₂ , r ₃
1 st char last name	r ₁₂	r ₃

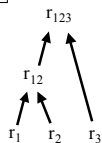


74

Iterative Blocking

Record	Name	Zip	Email
r ₁	John Doe	12345	jdoe@yahoo
r ₂	John Doe	94305	
r ₃	J. Doe	94305	jdoe@yahoo

Partition by	Block#1	Block#2
zip	r ₁ , r ₁₂	r ₂ , r ₃ , r ₁₂
1 st char last name	r ₁₂	r ₃

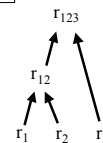


75

Iterative Blocking

Record	Name	Zip	Email
r ₁	John Doe	12345	jdoe@yahoo
r ₂	John Doe	94305	
r ₃	J. Doe	94305	jdoe@yahoo

Partition by	Block#1	Block#2
zip	r ₁₂	r ₁₂₃
1 st char last name	r ₁₂	r ₃

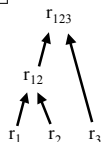


76

Iterative Blocking

Record	Name	Zip	Email
r ₁	John Doe	12345	jdoe@yahoo
r ₂	John Doe	94305	
r ₃	J. Doe	94305	jdoe@yahoo

Partition by	Block#1	Block#2
zip	r ₁₂₃	r ₁₂₃
1 st char last name	r ₁₂₃	r ₁₂₃

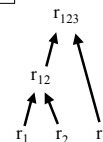


77

Iterative Blocking

Record	Name	Zip	Email
r ₁	John Doe	12345	jdoe@yahoo
r ₂	John Doe	94305	
r ₃	J. Doe	94305	jdoe@yahoo

Partition by	Block#1	Block#2
zip	r ₁₂₃	r ₁₂₃
1 st char last name	r ₁₂₃	r ₁₂₃



Iterative Blocking Solution: {r₁₂₃}

78

Overview

- Model
- Algorithms
- Experimental Results

79

Iterative Blocking Model

- Result is the fixed-point state of applying a “core” ER algorithm on blocks and re-distributing new records
- Can plug in any core ER algorithm that partitions records

80

In-memory Algorithm (Lego)

- Maintains “maximal records” for efficient block updates
- Maintains a queue of blocks to process

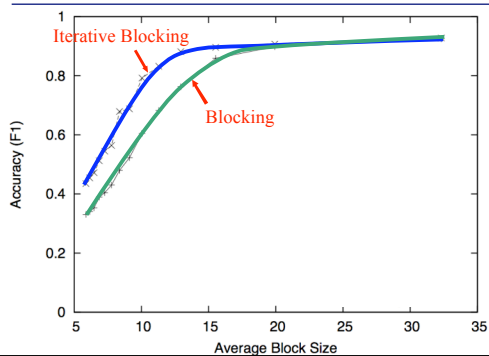
81

Disk-based Algorithm (Duplo)

- Reads into memory and processes “N blocks at a time” using segments
- Maintains a queue of segments to process
- Updates segments by scanning a merge log on disk
 - Uses timestamps to avoid a full scan for each segment read

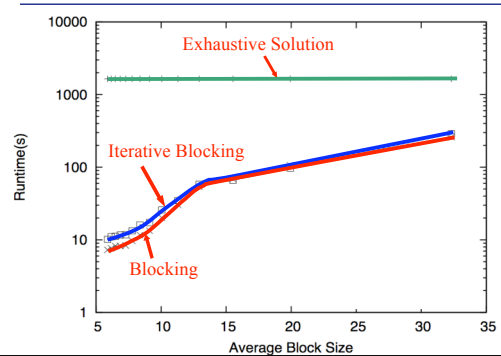
82

Accuracy

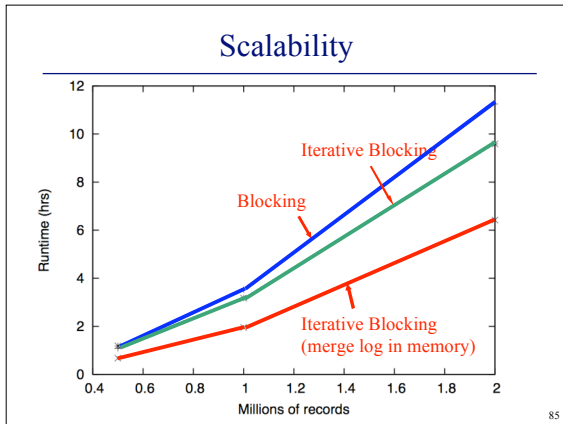


83

Performance



84



- ### Iterative Blocking Summary
- Proposed model & efficient algorithms (in-memory, disk) for iterative blocking
 - Showed that iterative blocking can be more accurate and scalable than simple blocking

- ### ER in the InfoLab
- Generic ER [VLDB J. 09]
 - Distributed ER [ICDCS 07]
 - Iterative Blocking [SIGMOD 09]
 - Negative Rules [VLDB J. 09]
 - Confidences [VLDB CleanDB 06]
 - Evolving Rules
 - Joint ER
 - ER Measures
 - ER in Probabilistic Databases
 - Privacy and Information Leakage