



Parallel Linkage



Hung-sik Kim¹ hungsik@psu.edu
Dongwon Lee^{1,2} dongwon@psu.edu

¹Computer Science and Engineering
²Information Sciences and Technology
 The Pennsylvania State University
 University Park, PA, 16802, USA


Record Linkage



- **Record Linkage:** Identify all matching records between two collections of records
 - Citations in digital libraries
 - Addresses in customer relationship management
 - Patient records in health information system
- Different from Database Join
 - Whole record vs. a few short attributes
 - Match vs. Match & Merge

CIKM 2007 / Parallel Linkage 2


Technical Landscape



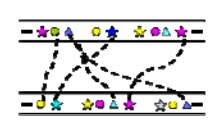
- Abundant research in many disciplines
- Also relevant to:
 - DB: approximate join, merge/purge
 - DL: citation matching, de-duplication, author name disambiguation
 - AI: identity matching
 - NLP: word sense disambiguation
 - IR: web query results clustering
 - LIS: name authority control

CIKM 2007 / Parallel Linkage 3

Related Work




- SERF project @ Stanford Univ.
 - Blackbox approach
 - Distributed and parallel record linkage
- Febrl project @ Australia National Univ.
 - Open source GUI based record linkage tool
- Data Linkage project @ Penn State Univ.
 - New applications
 - Group Linkage, Video Linkage
 - New aspects
 - Adaptability



<http://pike.psu.edu/linkage/>


CIKM 2007 / Parallel Linkage 4

Data Set




- Set characteristic
 - Clean set
 - A is clean iff $\forall r, s \in A, r \neq s$
 - Dirty set
 - A is dirty iff $\exists r, s \in A, r \approx s$
- Input data sets


Clean vs. Clean



Dirty vs. Clean




Dirty vs. Dirty







$O(|A||B|)$

CIKM 2007 / Parallel Linkage 5

Type of Matching & Merging



- $match(r, s)$ returns
 - $r \approx s$
 - $r \sqsubseteq s$  $\rightarrow merge(r, s) = s$
 - $r \supseteq s$  $\rightarrow merge(r, s) = r$
 - $r \equiv s$  $\rightarrow merge(r, s) = r$
 - $r \oplus s$  $r \cap s > \emptyset \rightarrow merge(r, s) = r \cup s$
 - $r \neq s \rightarrow$ no merge

CIKM 2007 / Parallel Linkage 6

Iterative Linkage Framework

- *match()*, *merge()* and iteration
 - a record is not similar to any records
 - a record is merged iteratively with a similar record

clean set clean set clean set clean set clean set

CIKM 2007 / Parallel Linkage 7

Proposed Algorithms

- Sequential Linkage: **six** variants
 - Clean-Clean: *S-CC*
 - Dirty-Clean: *S-DC*, *s-DC_{self}*
 - Dirty-Dirty: *s-DD1*, *s-DD2*, *s-DD3*
- Parallel Linkage: **five** variants
 - Clean-Clean: *p-CC*
 - Dirty-Clean: *p-DC*, *p-DC_{self}*
 - Dirty-Dirty: *p-DD1*, *p-DD2*

CIKM 2007 / Parallel Linkage 8

Sequential Linkage

- *s-DC* (Dirty vs. Clean)
 - *s-self* (one dirty set)
 - same as *s-DC*, but start with empty clean set

dirty set clean set dirty set clean set

Simple diagram

Iteration will stop when dirty set = \emptyset

CIKM 2007 / Parallel Linkage 9

Sequential Linkage

- *s-CC* (Clean vs. Clean)
 - iteration will stop when $C = \emptyset$
 - $s-DC_{self}(A, B) = s-CC(s-self(A), B)$
 - $s-DD1(A, B) = s-DC(A, s-self(B))$
 - $s-DD2(A, B) = s-CC(s-self(A), s-self(B))$
 - $s-DD3(A, B) = s-self(A \cup B)$

clean set A clean set B clean set A' clean set B' merge set C clean set A' union B merge set C

CIKM 2007 / Parallel Linkage 10

Parallel Linkage

- *p-CC*
 - put all clean sets into a buffer
 - partition one clean set, B, and send subsets
 - send the other clean set, A, to all processors

buffer: B_1, B_2, \dots, B_n, A

processor 1: A, B_1 processor n: A, B_n

processor 1: $A \cup B_1, C_1$ processor n: $A \cup B_n, C_n$

push C_1, C_2, \dots, C_n to a buffer

iteration will stop when the buffer is empty

CIKM 2007 / Parallel Linkage 11

Parallel Linkage

- *p-DC*: partition data sets and apply *s-DC*
- *p-DC_{self}*: partition data sets and apply *s-DC_{self}*
- *p-DD1*: partition data sets and apply *s-DD1*
- *p-DD2*: partition data sets and apply *s-DD2*
- **Similar implementation to *p-CC* (but much more complicated)**
 - Details omitted in the presentation

CIKM 2007 / Parallel Linkage 12

Experimental Setup

- Evaluation metrics
 - sequential linkage
 - number of comparisons (NC)
 - running time (RT)
 - parallel linkage
 - $speedup_{NC} / efficiency_{NC}$
 - $speedup_{RT} / efficiency_{RT}$
 - Data set
 - size : up to 5,000 records
 - Environment
 - distributed Matlab in PC cluster at Penn State
 - 1~32 nodes

CIKM 2007 / Parallel Linkage 13

Sequential Linkage

- Running time and Number of comparisons

CIKM 2007 / Parallel Linkage 14

p-CC

- Running time and Number of comparisons

CIKM 2007 / Parallel Linkage 15

p-CC

- Speedup and Scalability

CIKM 2007 / Parallel Linkage 16

Parallel Linkage

- Efficiency and Normalized RT/NC

CIKM 2007 / Parallel Linkage 17

Comparison against Others

- Efficiency_{RT} and Efficiency_{NC}

CIKM 2007 / Parallel Linkage 18

Conclusion



- Conclusion
 - Sequential linkage methods can be selected according to input data characteristics
 - Parallel linkage methods are shown to be efficient and scalable
 - Any distance function can be used in sequential or parallel linkage methods

Any Questions?



- Data/Codes are downloadable at:
 - <http://pike.psu.edu/download/cikm07/>



**Thank
You !**