



# Deep and Collective Entity Resolution in Parallel

Ting Deng   Wenfei Fan   Ping Lu   Xiaomeng Luo   Xiaoke Zhu   Wanhe An

Beihang University

University of Edinburgh

Shenzhen Institute of Computing Sciences



# Deep and Collective Entity resolution (ER)

- Identify whether two tuples in a dataset refer to the same real-world entity
- It has been a longstanding challenge to improve the accuracy of ER

*Deep and collective ER by unifying logic rules and ML models*

## *Deep ER (Recursion)*

Identify new matches  
by *making use of  
matches deduced earlier*

+

## *Collective ER*

Collectively correlate  
information *across  
multiple tables*

+

## *Unifying Logic and ML*

Benefit from both and  
improve the accuracy by  
*unifying the two methods*

# Matching Rules with mL (MRLs)

MRL  $\varphi: X \rightarrow I$

## Precondition $X$ :

A conjunction of *predicates* over a database schema  $\mathcal{R}$

$p ::= R(t) \mid t.A = c \mid t.A = s.B \mid \mathcal{M}(t[\bar{A}], s[\bar{B}])$

## Consequence $I$ :

An *id predicate*  $t.id = s.id$  or  
an *ML predicate*  $\mathcal{M}(t[\bar{A}], s[\bar{B}])$ .

- Extend matching dependencies by embedding with *well-trained ML classifiers* as ML predicates
- $t.id = s.id$ : the entities represented by  $t$  and  $s$  match, then refer to the same entity

## Deep and Collective ER

- Deduces a set  $\Gamma$  of *matches and validated ML predications* by applying MRLs in  $\Sigma$
- Modeled as an extension of the *chase* with MRLs, with the *Church-Rosser* property

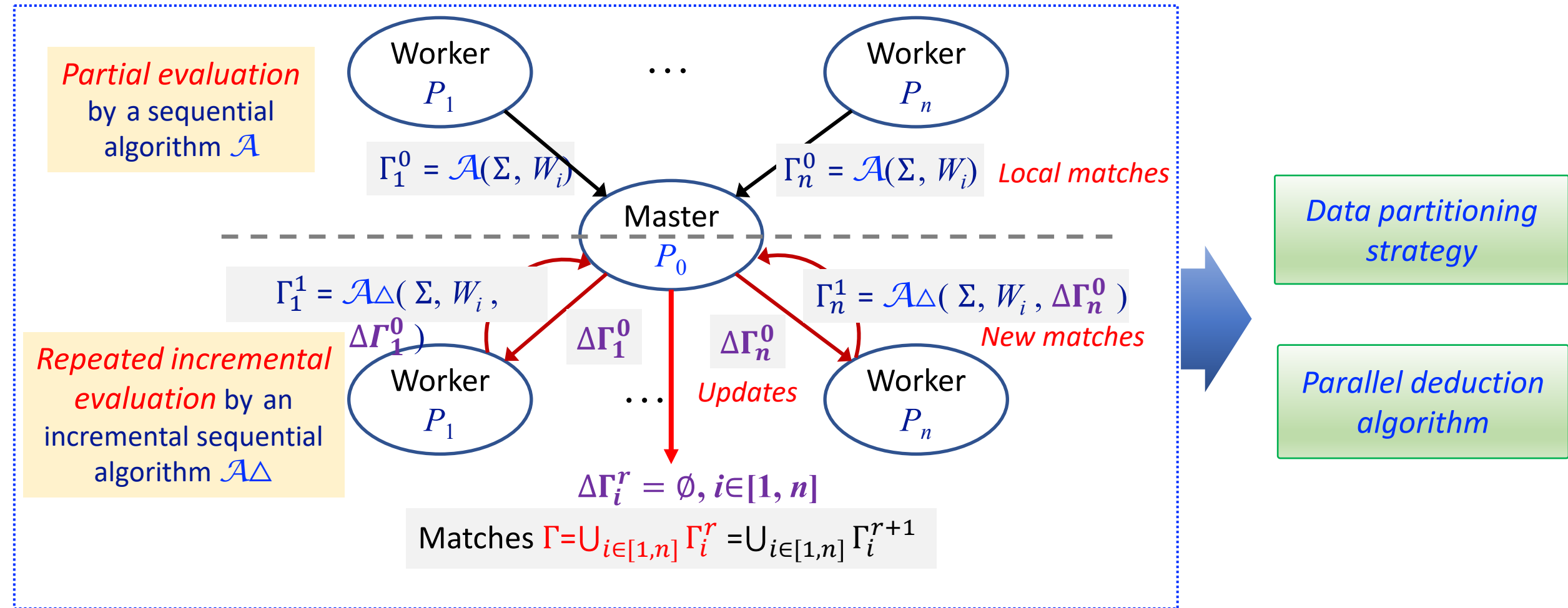
Decision problem	Number of relations	id predicates in the Prec.	
<b>Deep and Collective</b> ER	<i>Unbounded</i>	Yes	<i>NP-complete</i>
<b>Collective</b> ER (not Deep)	<i>Unbounded</i>	No	<i>NP-complete</i>
<b>Deep</b> ER (not collective)	Fixed	Yes	PTIME

*Parallel fixpoint  
model*



# Updated-driven Parallel fixpoint model under BSP model

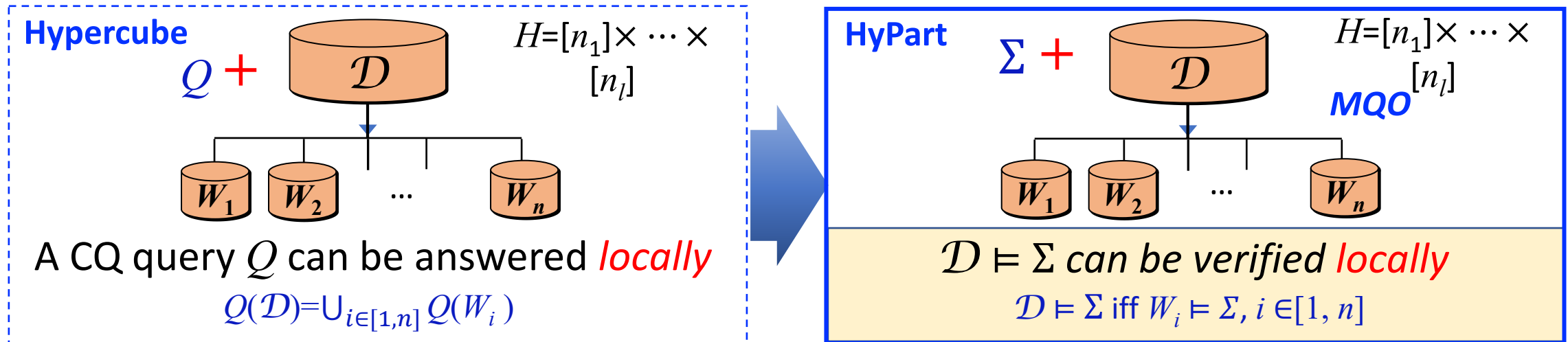
Partitioned dataset  $\mathcal{D} = (W_1, \dots, W_n)$ , distributed to workers  $P_1, \dots, P_n$



*Update-driven fixpoint computation: no communication of raw data*

# Data partitioning for parallel fixpoint computation

Extend Hypercube (HC) to handle a set of MRLs with *multiple query optimization (MQO)*



- Partition the data with the *minimum hash function computation (NP-complete)*

An heuristic partitioning algorithm **HyPart** / A strategy to assign the hash functions

Reduce the *computation* cost

*Reuse the hash function computations* as much as possible

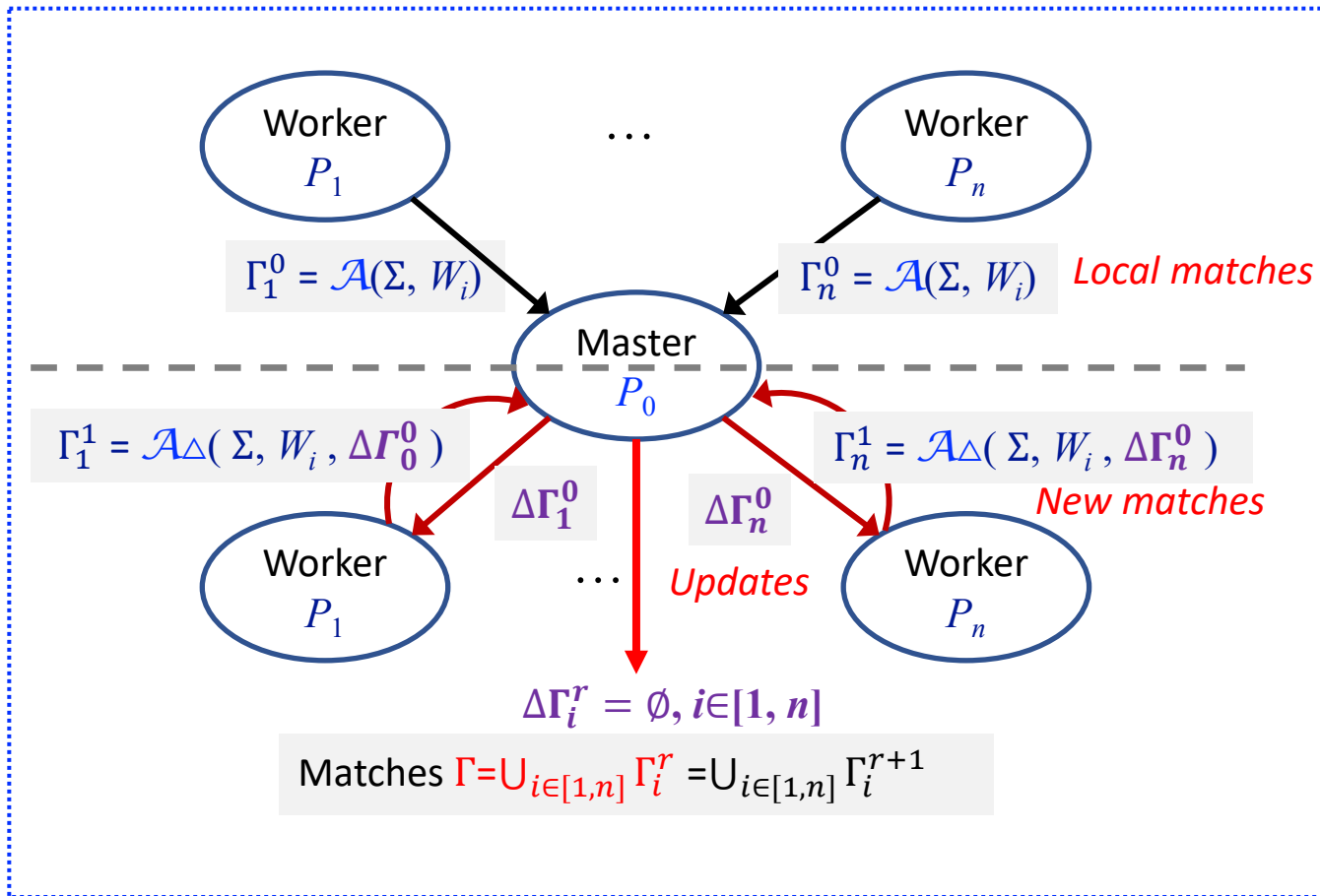
Reduce the *communication* cost

Send a tuple with the same hash functions *to the same worker* for different rules

Improve the performance

Using *MQO* to find the *common predicates* of MRLs

# A Parallel scalable deduction algorithm



*Partial evaluation by Deduce ( $\mathcal{A}$ )*

Deducing matching in  $\Gamma$  for each MRL and each of its valuations

*Repeated incremental evaluation by IncDeduce ( $\mathcal{A}_\Delta$ )*

Incrementally expand  $\Gamma$  by using an *update-driven* strategy

*Optimization to speed up the process*

Avoid to store all intermediate results and *reduce the repeated valuations*

*Parallelly scalable: the more processors are used, the faster the algorithms run*



# Experimental results

- Using *five real-life Datasets* and synthetic datasets, compared with *eight state-of-the-art* baselines

## Accuracy

- **23% and 38% more accurate** than ML and logical methods for ER, resp.
- outperforms deep ER and collective ER by 21% and 32%, resp.

## Efficiency

- **505s on datasets of 30M tuples** using 16 machines.
- faster than 7 out of 8 state-of-the-art ER baselines.

## Scalability

- **Parallely scalable** with the number  $n$  of processors used: 3.56 times faster when  $n$  increases from 4 to 32