The Interaction Between Schema Matching and Record Matching in Data Integration

Binbin Gu, Zhixu Li, Xiangliang Zhang, An Liu, Guanfeng Liu, Kai Zheng, Lei Zhao, and Xiaofang Zhou, Senior Member, IEEE

Abstract—Schema Matching (SM) and Record Matching (RM) are two necessary steps in integrating multiple relational tables of different schemas, where SM unifies the schemas and RM detects records referring to the same real-world entity. The two processes have been thoroughly studied separately, but few attention has been paid to the interaction of SM and RM. In this work, we find that, even alternating them in a simple manner. SM and RM can benefit from each other to reach a better integration performance (i.e., in terms of precision and recall). Therefore, combining SM and RM is a promising solution for improving data integration. To this end, we define novel matching rules for SM and RM, respectively, that is, every SM decision is made based on intermediate RM results, and vice versa, such that SM and RM can be performed alternately. The quality of integration is guaranteed by a Matching Likelihood Estimation model and the control of semantic drift, which prevent the effect of mismatch magnification. To reduce the computational cost, we design an index structure based on q-grams and a greedy search algorithm that can reduce around 90 percent overhead of the interaction. Extensive experiments on three data collections show that the combination and interaction between SM and RM significantly outperforms previous works that conduct SM and RM separately.

Index Terms-Data integration, schema matching, record matching

1 INTRODUCTION

UE to the data explosion in the big data era, the inconsistency between data sources becomes a critical issue in two dimensions: schema-level inconsistency and tuple-level inconsistency. As a result, merging data from multiple relational databases requires two necessary steps, namely Schema Matching (SM) and Record Matching (RM), in order to achieve a uniform and consistent data view. Here, SM unifies the schemas of different data sets; while RM finds pairs of linked records referring to the same entity.

There have been a host of works on SM or RM (see [31] or [14] for a survey). Briefly, the state-of-the-art SM method considers both the similarity (or semantic correlation) between the attribute names [20] and the similarity between the set of attribute values (or selected/sampled subsets of attribute values) under the two attributes [32]; while the most advanced RM methods inspect linguistic similarities and structural/relational similarities [3], [33] between key attribute values [4] or indicative non-key attribute values [38] when deciding the matching records among data sets, where key attribute is the one that can uniquely

- B. Gu, Z. Li, A. Liu, G. Liu, K. Zheng, and L. Zhao are with the School of Computer Science and Technology, Soochow University, Suzhou, Jiangsu 215006, China. E-mail: gu.binbin@hotmail.com, {zhixuli, anliu, gfliu, zhengkai, zhaol}@suda.edu.cn.
- X. Zhang is with the King Abdullah University of Science and Technology, Jeddah, Thuwal 23955-6900, Saudi Arabia. E-mail: xiangliang.zhang@kaust.edu.sa.
- X. Zhou is with the School of Information Technology and Electrical Engineering, The University of Queensland, Brisbane, QLD 4072, Australia, and the School of Computer Science and Technology, Soochow University, Suzhou, Jiangsu 215006, China. E-mail: zxf@itee.uq.edu.au.

determine a record in a relational table while all the others are non-key attributes. Recently, external domain knowledge and human interventions are also employed to improve the quality of SM [12] or RM [21].

All existing efforts, however, consider the two tasks independently, that is, they first perform SM, and then perform RM subsequently in only one run, which do not pay any attention on the possible interaction between SM and RM in the data integration process. This strategy inevitably gives us only one chance to make decisions, and deprives us further chances to update the links when more and more valuable information is collected from the other task. As a result, these SM and RM methods may easily make wrong decisions without further refined updates. Besides, there are two critical issues with existing works: (1) Existing instancebased SM methods rely heavily on the assumption that the distributions of attribute values under linked attributes should be similar to each other; otherwise it will suffer from low similarity between selected subsets of attribute values with the attributes should be linked, such as the situation in which only a small part of records are shared by the two data sets. (2) The RM linking results are greatly determined by SM linking results. As a result, both missed attributepairs and mistaken attribute-pairs would degrade the quality of RM linking results.

We study in this paper the interaction between SM and RM, by performing them alternately for data integration. To achieve this, novel matching rules are proposed: at each RM step, we identify a set of highly possible matching recordpairs based on the already linked attribute-pairs; Likewise, at each SM step, we identify a set of highly-possible matched attribute-pairs based on the already linked record-pairs. For instance, assume a start-up linked key attribute-pairs (Product, Product) between the two tables in Fig. 1a, at the first RM step, we may identify $(t_1 \leftrightarrow s_1)$ and $(t_2 \leftrightarrow s_2)$ as linked records as they share the same Product values. We

Manuscript received 21 Feb. 2016; revised 8 Sept. 2016; accepted 11 Sept. 2016. Date of publication 20 Sept. 2016; date of current version 5 Dec. 2016. Recommended for acceptance by K. Chang.

For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org, and reference the Digital Object Identifier below. Digital Object Identifier no. 10.1109/TKDE.2016.2611577



Fig. 1. Two example tables for integration (a) and the integration results with previous methods ((b) and (c)).



Fig. 2. Example interaction workflow of SM and RM for integrating tables in Fig. 1.

then identify (Weight, WT) as linked attribute-pair given that the two linked records share the same value under the two attributes. We repeat this process iteratively until no more attributes or records can be linked. Finally, we will have all the four attribute-pairs and six record-pairs be correctly linked as demonstrated in Fig. 2. By contrast, traditional methods perform SM and RM in only one run, which as a result introduce (Ex – Memory $\leftrightarrow \Rightarrow$ ROM) and $(t_4 \leftrightarrow \Rightarrow s_8)$ as wrong matches, and also miss pairs (Size; ***; Screen Size) $(t_3 \leftrightarrow s_3), (t_4 \leftrightarrow s_4), (t_5 \leftrightarrow s_5) \text{ and } (t_6 \leftrightarrow s_6) \text{ as}$ matched pairs with similarities and thresholds given in Figs. 1b and 1c, where the similarities between attribute values are measured by Levenshtein similarity. Instance-based SM methods in [13], [23] also use instances to facilitate SM. However, the instance-pairs are not from the results RM steps, and thus could be wrongly selected from mismatched attributes such as ROM and Memory in Fig. 1 due to the similar values they have.

Nevertheless, the interaction model raises two challenging issues: First, new linking decisions made at each SM (or RM) step based on intermediate RM (or SM) results should be *reliable*. Otherwise, they may lead to mistaken linking decisions in later stage. Second, the potential *semantic drift* in the interaction process should be controlled to prevent mismatch magnification in the subsequent iterations. Both of the two issues are crucial to the quality of the matching results.

To address the two problems, we design a probabilistic model to estimate the *Matching Likelihood* of each matchingrecord-pair. In particular, we first measure each individual attribute's ability to identify matching-record-pairs, and then estimate the likelihood of each matching-record-pair by jointly considering the identification ability of multiple attributes. The key difficulty lies on how to calculate the dependencies between attributes. A traditional model based on Inclusion-Exclusion principle [6] estimates the matching likelihood by calculating the dependencies between attributes comprehensively, but the computation cost grows exponentially with n (the number of attributes in a table). Another famous model Noisy-All [1], [26] completely neglects the dependencies between attributes for efficiency, with a sacrifice of estimation accuracy. To reach a balance between accuracy and efficiency, we propose a novel combination model which employs the logistic sigmoid function [5] to simplify the function of calculating the dependencies among the attributes into a linear one. Besides, to prevent from the semantic drift issue, we introduce different strategies to check the correctness of each matching-record-pair and matching-attribute-pair respectively. One checks the degree of deviation of every matching-record-pair from the other matching-record-pairs according to the unbiased variance [36], while the other employs cross-validation to use matching attribute-pairs to validate each other.

Computational cost is always an issue when comparing a large number of attribute value pairs in RM and SM. According to our analysis, without any optimizations, the computational complexity of the interaction algorithm can be as high as O(min(p,q)mn), where m and n are the numbers of records in the two tables for integration respectively, and p and q are the number of attributes in the two tables respectively. To reduce the high computational cost, we design an index structure based on q-grams [37] to index all possible matched record-pairs w.r.t. a single attribute. Potentially matchable record-pairs between the two tables are grouped into (possibly overlapped) blocks such that

matching-record-pairs are only identified within one block. We then propose a greedy algorithm selecting only one block at a time from all blocks, which brings the maximum benefit (i.e., linking the most matching-record-pairs or matching-attribute-pairs at the next step) with the minimum cost (i.e., comparing the least attribute values). After each step, the algorithm updates the indices and does the greedy block selection again until no more blocks left.

Our main contributions are summarized as follows:

- We first study the combination and interaction between SM and RM by performing them alternately when integrating multiple data sources. Novel matching rules are proposed as the foundation of the combination.
- 2) We design a probabilistic model based on the logistic sigmoid function to estimate the Matching Likelihood of a matching pair. Our model can be adopted to the situation with small overhead when many incidents are dependent, since a parameter acts on the logistic sigmoid function to smooth the dependency among attributes.
- 3) We propose two effective ways to check the correctness of each matching pair. One uses unbiased variance of the similarity between the attribute values pairs of the two records, while the other employs cross-validation to use all matching pairs to validate each other.
- 4) We greatly reduce the time complexity of the interaction algorithm by around 90 percent with a specialdesigned index structure and several optimization techniques proposed based on the index structure.

The rest of the paper is organized as follows: We give an overview of the interaction in Section 2. We discuss on the matching likelihood estimation scheme in Section 3, and then present how we control the integration quality in Section 4. The optimization on efficiency is given in Section 5. After reporting the experiments in Section 6, the related work is covered in Section 7. We conclude in Section 8.

2 INTERACTION OVERVIEW

Given two relational tables $T_1 = \{t_1, t_2, \dots, t_n\}$ under schema $S_1 = \{A_1, A_2, \dots, A_p\}$ and $T_2 = \{s_1, s_2, \dots, s_m\}$ under schema $S_2 = \{B_1, B_2, \dots, B_q\}$, where n, m, p, q are positive integers, t_i $(1 \le i \le n)$ denotes a record in T_1 , s_i $(1 \le i \le m)$ denotes a record in T_2 , A_i $(1 \le i \le p)$ denotes an attribute in T_1 , and B_i $(1 \leq i \leq q)$ denotes an attribute in T_2 , assume $S_1 \cap S_2 \neq \emptyset$, i.e., the two tables have common attributes, $T_1 \cap T_2 \neq \emptyset$, i.e., the two tables have records referring to the same real-world entity, the two fundamental tasks of Data Integration¹ is to perform Schema Matching between S_1 and S_2 , and Record *Matching* between T_1 and T_2 : While the objective of SM is to unify S_1 and S_2 by finding out all pairs of attributes (A_i, B_j) between S_1 and S_2 , where each attribute-pair refers to the same property of the records, or denoted as $(A_i \leftrightarrow B_j)$, the object of RM is to find out all pairs of records (t_i, s_j) referring to the same entity, or denoted as $(t_i \leftrightarrow s_i)$, between T_1 and T_2 . Here we assume that each attribute in S_1 matches with no more than one attribute in S_2 , and vice versa. Also, we assume

1. We present our interactive idea between two tables for ease of presentation, and it is extendable to the integration of multiple tables.

that each record in T_1 matches with no more than one record in T_2 , and vice versa. Besides, we do not consider the situation that an attribute in one table corresponds to the combination of several attributes in the other table.

The integration quality can be reflected in four dimensions, i.e., *Precision of SM*, *Recall of SM*, *Precision of RM*, and *Recall of RM*. In particular, the precision of SM is the percentage of correctly matched attribute-pairs among all matched attribute-pairs while the recall of SM is the percentage of correct matching-attribute-pairs among all matching-attribute-pairs that should be identified between the two tables. Similarly, the precision of RM is the percentage of correct matching-record-pairs among all matchingrecord-pairs while the recall of RM is the percentage of correct matching-record-pairs among all matchingrecord-pairs while the recall of RM is the percentage of correct matching-record-pairs among all matching-recordpairs that should be identified between the two tables.

To implement data integration, existing work does SM first and RM second in one run, but often fails to reach a satisfied integration quality (that could be achieved). In this paper, we work on the interaction between SM and RM in the process of performing them alternately, with the expectation that the interaction can provide us further chances to improve the integration quality. For easier presentation, we use *IntSRM* to denote the interaction process. In the following, we first use an example to demonstrate how we perform SM and RM alternately in Section 2.1, and then formally state the problems in Section 2.2.

2.1 Basic Interaction Scheme

The basic interaction process is described as follows: starting with seed linked attribute-pairs, we perform SM and RM in turn to detect linked attributes or linked records iteratively until no more links can be detected. While every RM step identifies more linked entities to help the next SM step find more not yet linked attribute-pairs, every SM step detects more linked attribute-pairs to benefit the followed RM in finding more not yet linked record-pairs.

Briefly, a RM matching is made based on temporary attribute-pairs that are already linked, and an SM matching is made based on temporary instance pairs that are already linked. More specifically, we describe the two rules below:

Rule 1. (*Linked-Attributes-based RM*). Given linked attributepairs $\{(A_{l1} \leftrightarrow B_{l1}), (A_{l2} \leftrightarrow B_{l2}), \dots, (A_{ll} \leftrightarrow B_{ll})\}$ between S_1 and S_2 in integrating T_1 with Schema S_1 and T_2 with schema S_2 , we say whether a record-pair $t \in T_1$ and $s \in T_2$ will be linked (temporarily) in the next RM step is determined jointly by the similarity between the pair $(t[A_i], s[B_i])$, and the ability of (A_i, B_i) in recognizing matching-record-pairs, where $i \in \{l1, l2, \dots, ll\}$.

Basically, the more attribute values the two records share under matching-attribute-pairs and the stronger the ability of these matched attribute-pairs in recognizing matchingrecord-pairs, the more likely that the two records can be a matching-record-pair referring to the same entity. For example in Fig. 1, given linked attributes (Weight $\leftrightarrow WT$) and (SIZE $\leftrightarrow Screen$). Assume the two attributes can effectively differentiate a record from the others, we may produce a record-pair ($t_4 \leftrightarrow s_3$), since it shares the same "Weight or WT" and "SIZE or Screen" values.

Rule 2. (*Linked-Records-based SM*). Given linked record-pairs $\{(t_{l1} \leftrightarrow s_{l1}), (t_{l2} \leftrightarrow s_{l2}), \ldots, (t_{ll} \leftrightarrow s_{ll})\}$ between T_1 and

 T_2 in integrating T_1 with Schema S_1 and T_2 with schema S_2 , we say whether an attribute-pair $A \in S_1$ and $B \in S_2$ will be linked (temporarily) in the next SM step is determined jointly by the similarity between each ($t_i[A]$, $s_i[B]$) pair, where $i \in \{l1, l2, ..., ll\}$.

The more record-pairs support the matching of the two attributes, the more likely that they should be matched pairs. If this situation is observed with several pairs of linked records between two tables, there will be a high matching likelihood that the two attributes actually refer to the same one. For example in the two tables in Fig. 1, Weight and WT might be linked since they share the same attribute values under several record-pairs like (t_1, s_1) and (t_2, s_2) .

Based on the two matching rules above, the interaction scheme can be simply represented by a sequence of attribute/ record-pair set. Let $\mathcal{P}_i^S = \{(A^{i,1} \leftrightarrow B^{i,1}), (A^{i,2} \leftrightarrow B^{i,2}), \ldots\}$ denote the attribute-pair set identified at the *i*th SM step, and $\mathcal{P}_i^R = \{(\mathbf{t}^{i,1} \leftrightarrow \mathbf{s}^{i,1}), (\mathbf{t}^{i,2} \leftrightarrow \mathbf{s}^{i,2}), \ldots\}$ denote the recordpair set identified at the *i*th RM step, an interaction scheme between T_1 and T_2 can be denoted as

$$Q = \langle \mathcal{P}_0^S, \mathcal{P}_1^R, \mathcal{P}_1^S, \mathcal{P}_2^R, \mathcal{P}_2^S, \dots, \mathcal{P}_k^R, \mathcal{P}_k^S, \dots \rangle,$$

where $\forall i \neq j, \mathcal{P}_i^S \cap \mathcal{P}_j^S = \mathcal{P}_i^R \cap \mathcal{P}_j^R = \emptyset$.

We now describe a basic interaction workflow between SM and RM with an example scenario in Fig. 2.

Example 1. Initially, we have $\mathcal{P}_0^S = \{(\text{Product} \leftrightarrow \text{Product})\}$, according to which we can match $(t_1 \leftrightarrow s_1)$ and $(t_2 \leftrightarrow s_2)$. Then, we find that the two matched records share the same values under (WT, Weight) and (SIZE, Screen). Thus, (WT $\leftrightarrow \text{Weight}$) and (SIZE $\leftrightarrow \text{Screen}$) can be our newly-linked attribute-pairs. Until now, we would have three attribute-pairs, according to which we can find a new record-pair $(t_4 \leftrightarrow s_3)$, given that the three matching-attribute-pairs support this record-pair. Next, since t_4 [CAMERA] equals to s_3 [BackCam] rather than s_3 [FontCam], we may have (CAMERA $\leftrightarrow \text{BackCam}$). We continue with RM and SM alternatively in this way to have $(t_5 \leftrightarrow s_4)$ and (ROM $\leftrightarrow \text{Memory}$).

2.2 Problems Statement

There are several crucial issues in the interaction workflow. First, the way of estimating the matching likelihood of an attribute-pair (or a record-pair) is the key factor to ensure the matching quality. As we mentioned in Rule 1, the matching likelihood between two records depends on two aspects, i.e., the number of linked attribute-pairs that support the matching, and the ability of the linked attributepair in recognizing matching-record-pairs. Therefore, the matching likelihood issue can be resolved by one sub-task of estimating the ability of the linked attribute-pairs in recognizing records referring to the same entity, and the other more challenging sub-task: how to combine the contributions from multiple attribute-pairs to the calculation of matching likelihood of the two concerned records. We will give our solution to this problem in Section 3.

Second, "semantic drift" problem should be controlled for preventing the mistake magnification from an SM (or RM) step to the following RM (or SM). The linking decisions made at each SM (or RM) step based on temporary RM (or SM) results should be validated. We will discuss the details of matching quality control in Section 4.

Last but not the least, the large overhead produced by comparing a large number of value pairs should be reduced. Our analysis shows that, without any optimizations, the computational complexity of the interactive algorithm can be as high as O(min(p,q)mn), where m and n are the number of records in the two tables for integration respectively, and p and q are the number of attributes in the two tables respectively. Section 5 will introduce how to reduce the computational cost.

3 MATCHING LIKELIHOOD ESTIMATION

We present the models for estimating the matching likelihood for record-pairs and attribute-pairs respectively.

3.1 Matching Likelihood Estimation for RM

We provide a way to estimate the ability of the linked attribute-pair in recognizing records referring to the same entity, and then discuss how to combine the contributions from multiple attributes to the matching likelihood of the two concerned records.

1) *IdC Score*. We call the ability of an attribute A in differentiating a record from the other records as the *Identification Confidence* of the attribute A, denoted as IdC(A). Basically, the IdC of an attribute can be learned from a large training data set with a probabilistic model, where the training data consisted of a set of labeled matched pairs denoted as Pos_T and a set of labelled unmatched pairs denoted as Neg_T . More specifically, we estimate the IdC score of an attribute A based on a labeled training set from its table as

$$IdC(A) = \frac{Pos_T(A)}{Pos_T(A) + Neg_T(A)},$$
(1)

where $Pos_T(A)$ is the number of record-pairs matched on attribute *A* among all labeled matched pairs in Pos_T , and $Neg_T(A)$ is the number of record-pairs matched on attribute *A* among all labeled unmatched pairs in Neg_T .

Given that two attributes, *A* from one table and *B* from the other, are matched, the IdC of the attribute-pair $A \iff B$ denoted by $IdC(A \iff B)$ can be estimated as

$$IdC(A \iff B) = \sqrt{\frac{[IdC(A)]^2 + [IdC(B)]^2}{2}}.$$
 (2)

2) Contributions Combination. Let $\{\mathbb{A} \leftrightarrow \mathbb{B}\} = \{(A_1 \leftrightarrow B_1), (A_2 \leftrightarrow B_2), \dots, (A_n \leftrightarrow B_n)\}$ denote the set of linked attribute-pairs, we discuss on how to estimate the matching likelihood of two records, say $t \in T_1$ and $s \in T_2$, based on: (1) the similarity of the values under linked attribute-pairs in s and t, denoted as $sim(s[A_i], t[B_i])$ ($1 \le i \le n$), and (2) the IdC score of every linked attribute-pair $(A_i \leftrightarrow B_i)$.

Two existing models are potentially usable for estimating the matching likelihood of two records, but have limitations discussed below.

(1) *Inclusion-Exclusion:* A classical way based on the inclusion-exclusion principle [6] calculates the matching likelihood of (t, s) as: $\mathcal{L}_{RM}(t, s) = \sum_{k=1}^{n} (-1)^{k+1} P(\Psi(k, n, t, s))$, where $\Psi(k, n, t, s)$ is a set containing all the *k*-size combinations generated from the set of linked attribute-pairs, and $P(\cdot)$ is the likelihood score

LRN The curve of our model is more smooth than the Noisy-All model and better estimates the matching likelihood as 1) it Our model 0.5 responds to a large span of x. while Noisy-All changes only largely within a small span of x; Noisy-All model 2) it takes small credits from attribute pairs who have negative overall contribution, while Noisy-All introduces will introduce negative values when adopting Eq (6), as shown by the dash line. x-Overall contribution of attribute pairs

Fig. 3. Comparison between models for illustration.

of a set as: $P(C_1 \cap \cdots \cap C_i) = \mu_{(C_1 \cap \cdots \cap C_i)} \cdot \prod_{j=1}^i IdC$ $(A_j \leftrightarrow B_j) \cdot sim(A_j, B_j) \cdot sim(t[A_j], s[B_j])$, where $C_j = A_j \leftrightarrow B_j$ and $\mu_{(C_1 \cap \cdots \cap C_i)}$ is the dependency factor of $\{C_1, \ldots, C_i\}$, which, however, needs to be estimated in advance. Although this model can accurately estimate the matching likelihood between two records by calculating the dependencies between attributes comprehensively, the computation cost grows exponentially with n.

(2) *Noisy-All:* The Noisy-All model [1], [26] is another popular model that calculates the matching likelihood as: $\mathcal{L}_{RM}(t,s) = 1 - \prod_{(A \leftrightarrow B) \in \{A \leftrightarrow B\}} sim(t[A], s[B]) \cdot (1 - IdC)$

 $(A \leftrightarrow B)$), which indicates that the likelihood estimation is simplified by assuming all attributes are independent. However, the dependency between attributes should not be neglected in real practice. In addition, the noisy-all model has the so called *Accumulative-Error* problem: the matching likelihood of record-pairs increases with the number of matched attributes. That is to say, two records can be wrongly estimated to have high matching likelihood even the similarity of values under their linked attribute-pairs is low, just because the number of matched attributes is high. For example, even $sim(t[A_i], s[B_i]) = 0.4$, $IdC(A_i \leftrightarrow B_i)) = 0.5$ $(1 \le i \le n)$ but n = 10, then $\mathcal{L}_{RM}(t, s) = 0.9$.

We propose a new model to take the advantages of both the two models but addressing their limitations. First, instead of the comprehensive way to calculate the dependencies among the attributes, we simplify the function into a linear one with the logistic sigmoid function [5] and then rely on only one parameter to control the influence among the attributes, i.e., we use the logistic sigmoid function to smooth the influence among the attributes for matching records in an explicit way. Besides, to overcome the Accumulative-Error problem, we define a contribution function and employ a logarithmic function to map the value from [0, 1] into $[0, +\infty)$.

To achieve this, we first assume that all the IdC of attributes are independent such that a linear model (similar to Noisy-All) can be used to calculate the matching likelihood, and then we compensate for the dependence between attributes in the model by introducing a damping factor. The matching likelihood of the record-pair (t, s) is then calculated as

$$\mathcal{L}_{RM}(t,s) = \frac{1}{1 + e^{-\lambda \cdot S(\mathbb{A} \leftrightarrow \mathbb{A}, t, s)}},$$
(3)

where λ is the damping factor to compensate for the dependence between attributes (the parameter can be tuned on a validation dataset), and $S(\mathbb{A} \iff \mathbb{B}, t, s)$ is the overall contribution score of the set of linked attribute-pairs $\{\mathbb{A} \iff \mathbb{B}\}$ to the matching of the record-pair (t, s), computed as

$$S(\mathbb{A} \iff \mathbb{B}, t, s) = \sum_{(A \iff B) \in \{\mathbb{A} \iff \mathbb{B}\}} \phi(A, B) \cdot ctr(t[A], s[B]), \quad (4)$$

where $\phi(A, B) \in [0, +\infty)$ employs a logarithmic function to map the value between 0 to 1 into the whole real axis as

$$\phi(A,B) = -ln(1 - \mathcal{L}_{SM}(A,B) \cdot IdC(A \iff B)), \quad (5)$$

where $\mathcal{L}_{SM}(A, B)$ is the matching likelihood of the two attributes A and B. Finally, ctr(t[A], s[B]) is the contribution of the similarity between two values t[A] and s[B]

$$ctr(t[A], s[B]) = \begin{cases} 0, & \text{if } t[A] = null \text{ or } s[B] = null \\ \frac{sim(t[A], s[B]) - \theta}{1 - \theta}, & \text{if } sim(t[A], s[B]) \ge \theta \\ \frac{sim(t[A], s[B]) - \theta}{\theta}, & \text{if } sim(t[A], s[B]) < \theta, \end{cases}$$
(6)

where sim(t[A], s[B]) is the similarity between t[A] and s[B]measured by string similarity function such as edit distance, and θ is an expert-defined tipping point to decide whether this value pair produces positive or negative contributions. The contribution defined in Eq. (6) resolves the issue of accumulative error when using sim(t[A], s[B]) directly as contribution. That is, a large but wrong $S(\mathbb{A} \iff \mathbb{B}, t, s)$ score can be obtained by accumulating small similarity sim(t[A], s[B]) of a large number of attributes $(A \iff B)$ in $\{\mathbb{A} \iff \mathbb{B}\}$.

To summarize, the proposed model has the following two advantages: (1) it can be adopted to the situation when the contributions of the linked attribute-pairs are not independent with a low overhead; (2) it solves the Accumulative-Error problem. As can be observed in the comparison of \mathcal{L}_{RM} curve between the Noisy-All model and our model in Fig. 3, the Noise-all \mathcal{L}_{RM} curve is abrupt and changes largely when x is small, while our \mathcal{L}_{RM} curve is smooth and responds to a large span of x. In addition, our model does not introduce negative values like Noisy-All model when adopting Eq. (6).

3.2 Matching Likelihood Estimation for SM

We now present how to estimate the matching likelihood between two attributes $A \in S_1$ and $B \in S_2$. Let $\{\mathbb{T}_t \rightsquigarrow \mathbb{T}_s\} = \{(t_1 \rightsquigarrow s_1), (t_2 \rightsquigarrow s_2), \ldots, (t_n \rightsquigarrow s_n)\}$ denote the set of linked record-pairs so far, we adopt the same model to calculate the matching likelihood between two attributes, where the damping factor is not required since the records are usually independent with each other. Thus we have

$$\mathcal{L}_{SM}(A,B) = \frac{1}{1 + e^{-S(\mathbb{T}_t \leadsto \mathbb{T}_s, A, B)}},\tag{7}$$

where $S(\mathbb{T}_t \leftrightarrow \mathbb{T}_s, A, B)$ is the overall contribution of the set of linked record-pairs $\mathbb{T}_t \leftrightarrow \mathbb{T}_s$ to the matching of the attribute-pair *A* and *B*, which can be computed as

$$S(\mathbb{T}_t \nleftrightarrow \mathbb{T}_s, A, B) = \alpha \cdot \sum_{(t \nleftrightarrow s) \in \{\mathbb{T}_t \nleftrightarrow \mathbb{T}_s\}} \varphi(t, s) \cdot ctr(t[A], s[B]),$$
(8)

where α is a parameter to control the contribution score about the related record-pairs, and $\varphi(t,s) = -ln(1 - \mathcal{L}_{RM}(t,s))$ is also a logarithmic function to map the value from [0, 1] into $[0, +\infty)$.

In the interactive process, our strategy is to keep a high precision with a strict constraint and improve the recall step by step. That is, at each RM or SM step, the \mathcal{L}_{RM} or \mathcal{L}_{SM} for a number of candidate record-pairs or attribute-pairs will be calculated, and only those satisfying predefined threshold will be taken as linked pairs for further interactions.

Example 2. Let (Product $\leftrightarrow Product$) be the seed attributepair to initiate the interaction between the two tables as depicted in Fig. 1. Assume the IdC(Product $\leftrightarrow Product$) is 0.97, and let 0.7 be the threshold to the matching likelihood of both SM and RM for illustration, $\theta = 0.1$ in Equation (6).

At the first RM step, we have $(t_1 \leftrightarrow s_1)$ and $(t_2 \leftrightarrow s_2)$ given that $\mathcal{L}_{RM}(t_1, s_1) = \frac{1}{1+e^{-3.5066}} = 0.97 > 0.7$ and $\mathcal{L}_{RM}(t_2, s_2) = \frac{1}{1+e^{-0.8723}} = 0.705 > 0.7$.

At the first SM step, we have (*Weight* $\leftrightarrow WT$) since they share the same two values in the two pairs of linked records such that the matching likelihood can be calculated as: $\mathcal{L}_{SM}(\texttt{Weight}, \texttt{WT}) = \frac{1}{1+e^{-[-ln(1-0.705)]}} = 0.991.$

4 QUALITY CONTROL

Although the interaction scheme introduced above tends to select more promising matching pairs for iterative interaction to keep a high matching precision, still, the linking decisions made during the interaction may involve errors, since: (1) the decisions made at each SM (or RM) step based on temporary RM (or SM) results should be selected for filtering out the unreliable ones; and (2) once a mistaken matching happens at an iteration, more mistaken matchings might be introduced in later iterations, i.e., the "semantic drift" problem happens.

In the following, we introduce how to control semantic drift, and how to iteratively update the linking pairs in the interaction for higher-quality linking results.

4.1 Semantic Drift Control

We validate the newly-linked records and newly-linked attributes separately to prevent semantic drift from happening. After each RM step, we identify "risky" record-pairs by checking the unbiased variance of the similarity between their value pairs under various attribute-pairs, while after each SM step, we identify "outlier" attribute-pairs by applying crossvalidation techniques to validate all the linked attributes.

1) Unbiased Variance Checking for "Risky" Records. Intuitively, two records having similar values under more attribute-pairs tend to have higher matching likelihood. However, relying on the matching likelihood only cannot effectively differentiate high-quality record-pairs from "risky" ones due to the ubiquitously existing errors, various formats and so on. To identify the risky pairs, we measure the degree of instability of each record-pair by calculating the variance of the similarity between their attribute values under various attribute-pairs. More specifically, for a record-pair ($t \leftrightarrow s$), we get the degree of instability under different attribute-pairs by calculating the Unbiased Variance [36] (short for UV) of the similarity between their value pairs under various attribute-pairs as

$$UV(t,s) = \frac{1}{m-1} \sum_{i=1}^{m} [sim(t[A_i], s[B_i]) - \overline{sim}(A_i, B_i)]^2, \quad (9)$$

where $\overline{sim}(A_i, B_i)$ is the average similarity of all the attribute values under the attribute-pair (A_i, B_i) and m is the number of linked attribute-pairs. We remove the record-pairs whose UV values are larger than a user-defined threshold, which is set as the average UV values of all the record-pairs to be checked in a data set by default. Note that $\overline{sim}(A, B)$ will change after removing some record-pairs, but we can obtain the final fixed record-pairs after several iterations. This process is similar to the k-means clustering [30] whose convergence property has been proved well. The complexity of the UV-checking method is $O(pq^2)$, where p is the number of linked attribute-pairs and q is the number of linked record-pairs to be checked.

This UV-checking is crucial to matching quality control, since it not only guarantees the high-quality record-pairs for next SM step, but also helps us identify different presentations of the same entity, for example "in" is an abbreviation of "inch" in the example in Fig. 1, since every pair of linked attribute values shares the same distance "ch".

2) Cross-Validation for Detecting "Outlier" Attributes. We adopt the cross-validation techniques to validate the linked attribute-pairs. Intuitively, if a pair of linked attributes does not consist with the other attribute-pairs, it is very likely a risky "outlier" that should be dropped.

Specifically, we denote the set of all attribute-pairs as \mathcal{P} , and partition \mathcal{P} into k disjoint subsets denoted as $\mathcal{P} = \{P_1, P_2, \dots, P_k\}$ and let the number of attribute-pairs in each subset in \mathcal{P} be $\frac{|\mathcal{P}|}{k}$ (We initialize k as the number of attribute pairs at one step. If $|\mathcal{P}|$ is too large, we can repartition them to reduce the computational cost). And we denote $\mathcal{P} - P_i$ the attribute-pair set which excludes P_i from \mathcal{P} . At each verification, we take $\mathcal{P} - P_i$ as a training set and P_i as a validation set. We exploit $\mathcal{P} - P_i$ to infer record-pairs, and then we compute the matching likelihood of attribute-pairs in P_i according to the inferred records pairs. Typically, we let the number of attribute-pairs in P_i be one (considering not too many attribute-pairs) and denote it as $(A \iff B)$. We adopt a linear loss function F in regression F(S, ((A, B), $(\eta)) = S_{(t \leftrightarrow s) \in \mathcal{R}}(t[A], s[B]) - \eta$, where \mathcal{R} is the record-pairs inferred by $\mathcal{P} - P_i$, *S* is a similarity function and η is a qualified threshold. And we define the error $\{0,1\}$ -loss in the judgement $F(S, ((A, B), \eta)) = 1$, if $S_{(t \leftrightarrow s) \in \mathcal{R}}(t[A], s[B]) < \eta$. Then we calculate the errors (short for *Er*) of the attributepair (A, B) as follows:

$$Er(A,B) = \frac{1}{|\mathcal{R}|} \sum_{(t \leftrightarrow s) \in \mathcal{R}} F(S, ((A,B),\eta)).$$
(10)

We repeat this verification process k times with different validation sets and then we obtain the Er values of each attribute-pair. We then drop the attribute-pairs whose Er is lower than a predefined threshold.

4.2 Iterative Updating and Adjusting

As the interaction proceeds, more and more attribute-pairs and record-pairs are linked. After each iteration, the matching likelihood we calculated in previous iterations need to be updated, according to which we also need to adjust the attribute and record-pairs that are already linked.

The relationship between the record-pairs and attributepairs is mutually reinforced.We can use a bi-graph to illustrate the relationship between the record-pairs and attribute-pairs, where the weight on an edge means the contribution score of a record-pair or an attribute-pair. Typically, the weight on an edge which points to an attributepair (A, B) from a record-pair (t, s) is $\varphi(t, s) \cdot ctr(t[A], s[B])$, and the weight on an edge which points to a record-pair (t, s) from an attribute-pair (A, B) is $\phi(A, B) \cdot ctr(t[A], s[B])$. Then the matching likelihood of a record-pair is to apply the function $f(x) = \frac{1}{1+e^{-\lambda \cdot x}}$ to the summation of all the weights on the edges pointed to the record-pair itself. And this calculation method can be also adapted to an attributepair. Denote the matching likelihood of record and attribute-pairs respectively as a vector \vec{r} and \vec{a} . We compute \vec{r} and \vec{a} alternatively until they all reach a stable state. We prove that the vectors \vec{r} and \vec{a} will converge to a constant vector. Formally, we have the following conclusion:

Theorem 1. *The iterative algorithm is convergent. That is, the matched record and attribute-pairs will be uniquely determined finally.*

Proof. Please see the appendix.²

4.3 Algorithm Analysis and Other Issues

One issue of the algorithm is how to set the thresholds for the matching likelihood of record-pairs and attribute-pairs respectively. Although we set the thresholds empirically for different data sets, the experiment results show that the thresholds can be set to values in a quite large range where the performance of our algorithm has no significant variation. The reason is that the initial threshold setting will not affect the number of matched record or attribute-pairs at the later steps. Specifically, if the threshold is too high, we just need more iterations to find more matched record or attribute-pairs. If the threshold is too low, the quality control process will gradually amend them as the interaction algorithm goes. Generally, the thresholds should be larger than 0.5, which means one record or attribute-pair is more likely to be a correct one than a wrong one. In our experiments, we set both of them to 0.6. Also, the threshold setting for Er is similar to the threshold setting to matching likelihood.

The time complexity of the algorithm is mainly decided by both the time complexity of SM steps and that of RM steps. The cost of SM mainly depends on the comparison times between two subsets of instances that are used for schema matching, while the cost of RM mainly depends on the number of records and the size of schemas in the two tables. Formally, we analyze the time complexity of our algorithm as follows:

Theorem 2. Given tables $T_1 = \{t_1, t_2, ..., t_n\}$ under schema $S_1 = \{A_1, A_2, ..., A_p\}$, and $T_2 = \{s_1, s_2, ..., s_m\}$ under

2. Due to the limitation of space, we put our appendix online which can be accessed through http://ada.suda.edu.cn/Uploads/File/201512/04/1449212591654/proof.pdf

schema $S_2 = \{B_1, B_2, ..., B_q\}$ for integration. The upperbound of the time complexity that an interaction scheme for IntSRM can reach is: O(min(p,q)mn), where p and q are the number of records in T_1 and T_2 respectively, m and n are the number of attributes in S_1 and S_2 respectively.

Proof. Please see the appendix.

П

5 OPTIMIZATIONS ON EFFICIENCY

As can be seen from the proof of Theorem 2, the efficiency bottleneck of the interaction usually lies on the RM step since there are usually a lot more records than attributes in the table. To minimize the number of record-pairs for comparison, some state-of-the-art indexing techniques [8] have been proposed for scalable record linkage and successfully applied on RM based on the key attribute. In this paper, we extend the q-gram index [9], [22], [37] to multiple pairs of attributes scenario, and split potential matched record-pairs between the two tables into (possibly overlapped) blocks so that matching-record-pairs are only identified within every block.

5.1 Indices for SM and RM Interaction

Before introducing how the indices are built on linked attribute-pairs, some definitions and lemmas are given first:

Definition 1. Given a string s, a set of q-grams (q is a constant to denote the length of each gram) can be generated from s as: $Gms(s,q) = \{gm_1, gm_2, \dots, gm_{|s|-q+1}\}$, where gm_i consists of the characters from i to (i + q - 1) by their natural order in s. Then an ℓ -length consecutive q-gram sequence (or (ℓ, q) -seq for short) of s can be defined as a string that consisted of a sequence of q-grams consecutively in their natural order in Gms(s,q).

Lemma 1. Let $F(s_1, s_2) = \frac{|Gms(s_1,q) \bigcap Gms(s_2,q)|}{|Gms(s_1,q) \bigcup Gms(s_2,q)| + \varepsilon'}$ where ε gets

very close to 1. *Given a q-gram overlap threshold* ω *, if* $F(s_1, s_2) \geq \omega$ *, then they should share at least one* (ℓ, q) *-seq, where* $\ell \geq q \cdot \lfloor (max(|s_1|, |s_2|) - q + 1) \cdot \omega \rfloor$.

Proof. Please see the appendix.

1) Index Building on a Single Pair of Attributes. Based on the definition and the lemma above, the index building process can be described as follows: Given a pair of linked attributes, for every distinct attribute value s under the linked attributes in either table, we generate all (ℓ, q) -seq from this value, where $\ell \ge q \cdot \lfloor (|s| - q + 1) \cdot \omega \rfloor$, where *q* is a constant to define the length of grams. Based on Lemma 1, each (ℓ, q) -seq of an attribute value *s* will be taken as a *key index* value for s, according to which s will be indexed into a corresponding block. For instance, assume a product value "huawei", whose bi-gram list is ['hu', 'ua', 'aw', 'we', 'ei']. Let $\omega = 0.8$, then $\ell \ge q \cdot |(|s| - q + 1) \cdot \omega| = 8$, thus we will generate six $(\ell, 2)$ -seq from "huawei", i.e., 'huuaawweei', 'uaawweei', 'huawweei', 'huuaweei', 'huuaawei', 'huuaawwe' as the key index values for "huawei". As a result, "huawei" will be put into the six blocks corresponding to the six key index values.

2) Dynamic Indices Building Between the Two Tables. It is a dynamic process to build the index on multiple pairs of linked attributes in the interaction process. Initially, we only build the index under the seed attribute-pair (such as (Product $\leftrightarrow Product$)). When more linked attributes are identified at each SM step, we build an index under each of these linked attributes, as long as its IdC score is higher than a predefined threshold. For easier presentation, we call the index we build on the two databases for integration as a *Qgram-based Multiple-Line Indices for the Interaction between RM and SM* (or *MLineIndex* for short).

Lemma 2 describes the relationship between the threshold ω we mentioned above and the edit similarity threshold between two attribute values.

Lemma 2. Let ω denote the threshold that works for controlling the generation of key value indexes in building the MLineIndex, if two attribute values s_1 and s_2 are assigned into one block, then their edit similarity must be no less than $\frac{\lfloor (max(|s_1|,|s_2|)-q+1)\cdot\omega\rfloor+q-1}{max(|s_1|,|s_2|)}$, and if s_1 and s_2 are not in the same block, their edit similarity must be no larger than $1 - \frac{2-q+max(|s_1|,|s_2|)-\lfloor(max(|s_1|,|s_2|)-q+1)\cdot\omega\rfloor}{q\cdot max(|s_1|,|s_2|)-(max(|s_1|,|s_2|)-q+1)\cdot\omega]}$.

Proof. Please see the appendix.

5.2 Greedy RM Based on MLineIndex

We now describe a greedy RM algorithm with the MLineIndex at a particular RM step. Assume we already have a set of linked attributes, and every potential matched record-pairs w.r.t. a linked attribute-pair is put into a block and indexed under this linked attribute-pair. For the sake of processing the interaction between SM and RM with the minimum RM comparison times, each RM step only greedily selects one particular block of records for RM comparison from all the blocks indexed by the MLineIndex, which should satisfy the following two conditions: (1) it requires the least RM comparison times; (2) it has a high probability to generate matching record-pairs. More specifically, given a block *Block* = ({*LR*}, {*RR*}), where *LR* is the set of records from one table, and *RR* is the set of records from the other table, we estimate the priority of doing RM to this block as

$$priority(Block) = \frac{IdC(AttrPair_{Block})}{Max(|LR|, |RR|)},$$
(11)

where $AttrPair_{Block}$ is the attribute-pair under which the block is indexed, and $Max(|LR|, |RR|) = \frac{|LR| \times |RR|}{Min(|LR|, |RR|)}$ denotes the average comparison times needed for generating a matching record-pair from Block, since the block contains at most Min(|LR|, |RR|) record-pairs. When a block $Block = (\{LR\}, \{RR\})$ is selected for RM comparison, for every pair of records between LR and RR, we calculate the similarity between their attribute values under each linked attribute-pairs respectively, and then get their RM matching likelihood according to Eq. (3). We set a maximum number of record-pairs to be used for SM (1,000 in our experiments), such that the time consumption of SM steps will not be deferred too much by RM steps.

We observe that the block with higher priority often produces correct record-pairs based on MLineIndex, although we have not considered the similarity of attribute values in Eq. (11). The reason is that the high priority of a block often means there are less records in the block, but they often have quite a lot same q-grams, i.e., their similarities are often very large.

Another advantage of this strategy is that we can rapidly identify the matched record-pairs, since many attributepairs have been identified. Such that, for the rest records,



Fig. 4. The MLineIndex built on the two example tables in Fig. 1.

they can reduce many unnecessary comparisons due to the less and less number of records.

How to set the threshold ω_i for each attribute-pair A_i is the key problem for the indices-based greedy RM algorithm. An optimal setting is demanded to make a trade-off between the precision of matching results and the efficiency. In the following, we discuss the setting of ω_i for maximizing the precision of matching results while reducing the time complexity as much as possible.

5.2.1 Setting the Parameter ω_i

Computing time cost is mainly affected by the number of blocks and the number of record-pairs in blocks. Generally, to reduce the number of blocks and the record-pairs in blocks, it is preferred to have a relatively high threshold ω to those attribute-pairs with a relatively low *IdC*. As the number of q-grams generated from an attribute values grows exponentially as the threshold ω decreases [8], we employ a convex decreasing function to set the threshold for an attribute-pair $A \iff B$ according to its *IdC*. In particular, let $\{A_1 \iff B_1, A_2 \iff B_2, \ldots, A_m \iff B_m\}$ denote the set of matched attribute-pairs sorted in a non-decreasing order of their *IdC* scores, and let ω_i $(1 \le i \le m)$ denote the q-gram overlap threshold for the attribute-pair $A_i \iff B_i$, we choose ω_1 as a base value (it will be set automatically in later discussions), and then obtain ω_i with the following equation:

$$\omega_i = e^{\frac{IdC(A_i \leftrightarrow B_i) \cdot ln\omega_1}{IdC(A_1 \leftrightarrow B_1)}} \tag{12}$$

and we will show how to set ω_1 below.

According to Lemma 2, if two records t and s are not in a same block pair under the index of $A_i \leftrightarrow B_i$, the similarity between their attribute values under $A_i \leftrightarrow B_i$ has a upper bound denoted by sim_i^u . Thus we can rewrite sim_i^u as the following equation combining with Eq. (12)

$$sim_{i}^{u} = 1 - \frac{2 - q + M_{i} - \lfloor (M_{i} - q + 1) \cdot \omega_{1}^{\frac{IdC(A_{i} \sim -B_{i})}{IdC(A_{1} \sim -B_{1})} \rfloor}{q \cdot M_{i}}, \quad (13)$$

where M_i is the maximal length of values under $A_i \leftrightarrow B_i$. Thus if one record-pair (t, s) is not in any block pair of all the indices, we can obtain its upper bound of matching like-lihood as follows:

$$\mathcal{L}_{RM}^{U}(t,s) = \frac{1}{1 + e^{-\lambda \cdot \sum_{i=1}^{m} -\ln(1 - IdC(A_i \leftrightarrow B_i)) \cdot ctr(t[A_i], s[B_i])}}.$$
 (14)

Recall Eq. (6), since $sim(t[A_i], s[B_i])$ is larger than θ , $ctr(t[A_i], s[B_i])$ here can be substituted with $\frac{sim_i^u - \theta}{1 - \theta}$. Given a quality threshold τ_p^R to matched record-pairs, and let $\mathcal{L}_{RM}^U(t, s) = \tau_p^R$, we can derive the value of ω_1 by replacing sim_i^u with Eq. (12), and ω_1 will satisfy the following equation:

$$\frac{1}{1 + e^{-\gamma \cdot \sum_{i=1}^{m} -\ln(1 - IdC(A_i \nleftrightarrow B_i)) \cdot (sim_i^u - \theta)/(1 - \theta)}} = \tau_p^R.$$
 (15)

Eq. (15) is an equation with only one unknown parameter ω_1 . Therefore, ω_1 can be derived by analysing Eq. (15). And all the other thresholds about ω_i can be derived according to Eq. (12).

Theorem 3. By setting w_1 which satisfies Eq. (15) and w_i $(1 \le i \le m)$ according to Eq. (12), all possible matched pairs (whose matching likelihood is larger than τ_p^R) can be covered by the minimum number of blocks under linked attribute-pairs.

Proof. Please see the appendix.

5.2.2 Bounding for RM Step

Although setting thresholds can prune a large percent of unmatched record-pairs for comparison, there are still many unmatched record-pairs waiting to be filtered in the blocks. Here we employ a bounding-based strategy to further identify unmatched record-pairs from matched ones.

For a record-pair (t, s) in a block under a set of attributepairs $\{A_i \nleftrightarrow B_i\}$, where $1 \le i \le m$, it satisfies $\frac{\lfloor (max(|s_1|,|s_2|)-q+1)\cdot \omega \rfloor + q-1}{max(|s_1|,|s_2|)} \le sim(t[A_i], s[B_i]) \le 1$ according to Lemma 2. Then the lower bound of the matching likelihood $\mathcal{L}_{RM}^L(t,s)$ between t and s can be calculated with the lower bound of $sim(t[A_i], s[B_i])$, while the upper bound of the matching likelihood between t and s is

$$\mathcal{L}_{RM}^{U}(t,s) = \frac{1}{1 + e^{-\lambda \cdot \sum_{i=1}^{m} -ln(1 - IdC(A_i \leftrightarrow B_i))}}.$$
 (16)

If $\mathcal{L}_{RM}^{L}(t,s) > \tau_{p}^{R}$, then (t,s) will be a candidate record-pair, and if $\mathcal{L}_{RM}^{U}(t,s) < \tau_{p}^{R}$, then (t,s) will be pruned directly. In order to tighten the bounds as early as possible, we compute the similarity of the attribute values under the attribute-pair with higher *IdC* in priority, and then the upper-bound and lower-bound of the likelihood can be updated correspondingly.

After all, we will analyze the time complexity of this greedy and bounded interaction algorithm based on the MLineIndex in the appendix.

6 **EXPERIMENTS**

6.1 Data Sets and Metrics

We conduct experiments on two real and one synthetic data sets:

• *Mobile:* We collect cellphones on sale from Tmall³ and PConline⁴ respectively. The Tmall table contains

40k tuples under 53 attributes, while the PConline table contains 56k tuples under 46 attributes. The two tables share 3.8k records and 38 common attributes including *Release Date, Operation System, RAM, Screen Size, Type* etc.

- *Camera:* We collect digital cameras on sale from Yesky⁵ and PConline respectively. The Yesky table contains 25k tuples under 50 attributes, while the PConline table contains 34k tuples under 44 attributes. The two tables share 25k records and 31 common attributes including *Type, Pixels, Panel, Wifi, Manufacturer* etc.
- *Synthetic:* We also generate two synthetic tables sharing 100k tuples and 60 common attributes and use certain rules to let the distribution of data close to real data sets. For instance, the similarities between attribute values in linked record-pairs are uniformly distributed between 0 and 1. Note that since there are some missing attribute values in the two real data sets, we also generate a random number of missing values under each non-key attribute for the synthetic data set.

Metrics. We evaluate the effectiveness of the integration methods in four dimensions, i.e., *Precision of SM, Recall of SM, Precision of RM,* and *Recall of RM.* Also, the F_1 -score of SM and RM are also concerned. We use *Time Cost* to evaluate the efficiency of a method.

6.2 Integration Quality Comparison

We compare the integration quality of our method (short for *IntSRM*) with several state-of-the-art methods on the three data sets. For a fair comparison, for each method testing on every dataset, we tuned its setting to make the method reach the best performance on that specific dataset.

- a) Name-based SM (*NBSM*): This SM method uses the edit similarity between the attribute names [20] for SM.
- b) Value-based SM (*VBSM*): This SM method uses the overlap between the selected subset of attribute values under the two attributes [19] for SM, where each selected subset contains the top-k highest frequency attribute values under the attribute.
- c) Linkage-Points-based SM (*LPSM*): This is a state-ofthe-art instance-based SM method proposed in [23], which treats the matching function as a black-box and uses specific measures to have reliable SM results from the overlapping instances.
- d) Key-based RM (*Key*): This RM method inspects the string similarities between key attribute values [4] only for RM, which also uses q-gram together with inverted index [34], as well as prefix-based pruning [35] and batch-based matching [7], for improving efficiency.
- e) Key+Non-key RM (*NokeaRM*): This is a state-of-the-art RM methods using both key and non-key attributes for RM [11], [38]. Briefly, it builds a probabilistic rulebased decision tree based on all attributes according to the ability of each attribute in identifying matching





Fig. 5. Comparing the F_1 -score of RM and SM methods, respectively, on the three data sets.



Fig. 6. Comparing the precision and recall on three data sets.

records and the ability in identifying un-matching records, and then relys on this decision tree to make RM decisions. To further improve its efficiency, q-gram-based blocking techniques are used: we generate q-grams for each value, and only those satisfying the minimum q-gram overlap will be compared.

In the following, we first compare the integration quality of RM with previous methods, and then that of SM with previous methods. Since the overlap ratio of the records between two tables has a great influence on the integration quality, we conduct our comparison experiments at various overlap ratios (from 10 to 90 percent) on the three data sets.

(1) F_1 Comparison for RM. As demonstrated in Figs. 5a, 5b, and 5c, pervious RM methods work poorly (with F_1 -score around 0.5-0.6 for NokeaRM and F_1 around 0.2-0.4 for Key) when the overlap ratio is low (such as 10, 30 percent), while IntSRM can reach F_1 -score as high as 0.6-0.8 on all the three data collections. As the overlap ratio increases, the integration quality of all methods increases gradually. But always, our method IntSRM reaches 20 percent higher F_1 than the other methods.

(2) F_1 Comparison for SM. Similar comparison results can be observed for SM. As demonstrated in Figs. 5d, 5e, and 5f, pervious SM methods work poorly (with F_1 -score less than 0.6) when the overlap ratio is low (such as 10, 30 percent), while IntSRM can reach 0.7-0.8 on all the three data collections. As the overlap ratio increases, the integration quality of all methods increases gradually. But always, our method IntSRM reaches about 15 percent higher F_1 -score than the other methods.

(3) *PR* (*Precision & Recall*) *Comparison for SM and RM Respectively.* For more comprehensive comparison, we also draw the Precision-recall graphs for RM and SM comparison by setting the overlap ratio 70 percent on all the three data collections. As shown in Fig. 6, basically, our method IntSRM can always reach the best performance on either precision or recall over the three data collections.

6.3 Efficiency Comparison

We compare the efficiency of our method with previous methods at the setting of the overlap ratio 70 percent. For fair comparison, we also compare the time cost of SM or RM with previous SM or RM methods separately. For IntSRM, the time cost of SM includes not only the time cost of all SM steps, but also the time cost of all the RM steps that processed before the last SM step. The time cost of RM for IntSRM includes the time cost of all SM and RM steps.

(1) *Time Cost for RM.* As described in Figs. 7a, 7b, and 7c, the Key method uses much less time than NokeaRM and IntSRM since the Key method uses the key attribute only for RM. Compared with NokeaRM, IntSRM uses a bit more time since NokeaRM employs a special decision tree to help prune a large percentage of record-pairs for comparison, which on the one hand, greatly improves the efficiency, but on the other hand, hurts the precision and recall as reported



Fig. 7. Comparing the time cost with other methods on three data sets.



Fig. 8. Quality improvement on "Mobile" and "Camera" data sets and loss of quality on "Mobile" data set.



Fig. 9. Quality improvement with interaction (on Mobile and Camera) and parameter setting test (on syn).

in the experiments in the last section. Overall, the time cost spend in RM by IntSRM is acceptable as we greatly improve the precision and recall of the integration.

(2) *Time Cost for SM*. As described in Figs. 7e, 7f, and 7g, VBSM uses the least time than both LPSM and IntSRM. The time cost of IntSRM is more than LPSM since the IntSRM spends more time on finding record-pairs under more attribute-pairs than LPSM. But the time cost of IntSRM is acceptable in practice since we can greatly improve the precision and recall.

6.4 Quality Improvement Evaluation

We evaluate the effectiveness of our proposed techniques for improving the quality of the integration results. For differentiation, we call the interaction without any quality control techniques as the *Baseline*, the one with validating newlinked records as *UV-Check*, the one with validating newlinked attributes as *Cross-Validate*, and the one with both the two techniques as *UV+Cross*.

(1) Semantic Drift Control. As demonstrated in Figs. 8a, 8b, 8c, and 8d, the two techniques have different performance

on the two data collections, but the combination of them apparently improves the integration quality of both SM and RM. Specifically, UV-Check tends to improve the precision of the Baseline by around 10 percent for SM without hurting the recall, while Cross-Validate tends to improve the precision of the Baseline by around 10-15 percent but may hurt the recall of SM. This is because Cross-Validate uses a strict rule in deciding unqualified matching pairs. We set the threshold of *Er* value as 80 percent here. Overall, the combination of the two always improves the precision and recall of both SM and RM.

(2) *Iterative Updating*. Fig. 9b shows the quality of RM and SM can make a further improvement and they can hold steady with a satisfied result as the iteration goes.

6.5 Efficiency Improvement Evaluation

We evaluate the effectiveness of our proposed techniques on improving the efficiency of our interaction algorithm. We use *Baseline* to denote the algorithm without any optimization on efficiency, and *Greedy* to denote our interaction algorithm based on the MLineIndex we build.



Fig. 10. Effect of missing values to integration quality.

(1) *Efficiency Improvement*. As shown in Fig. 7d, we can find that our Greedy method actually can save almost 90 percent time cost of the baseline which proves the effective-ness of the proposed techniques.

(2) *Side-Effect: Loss in Quality.* As a side-effect of the efficiency improvement, there is a little decrease on both precision and recall. As shown in Figs. 8e and 8f, there is less than 5 percent decrease on precision and 2-5 percent decrease on recall, which is an acceptable price to pay for the reduction of 90 percent of the time cost.

6.6 Effect of Missing Values

We also conduct experiments on evaluating the effect of missing values to the performance of the proposed approaches and the other methods by randomly removing some non-key values from the table. As can be observed in Figs. 10a and 10b, as the missing ratio increases from 0 to 60 percent, the integration quality (including the F_1 -score of SM and RM) of all approaches using non-key attribute values decreases. For SM, LPSM decreases the most from more than 0.80 to about 0.57, while our approach IntSRM can still reach about 0.88 when the missing ratio becomes 60 percent. For RM, NokeaRM is also robust and only decreases from about 0.7 to about 0.63, while our approach IntSRM decreases from 0.89 to 0.75. Generally, our method always reaches 10 percent higher F_1 -score than the other methods.

6.7 Parameter Setting

We evaluate the effect of the three key parameters in our algorithm to the quality of data integration, which are: (1) the threshold to the SM matching likelihood; (2) the threshold to the RM matching likelihood; and (3) the threshold of Er value of Cross-Validate. To avoid potential biases, we fix the other two when we evaluate one of the three parameters. As can be seen from Figs. 9c and 9d, the F_1 -score of SM or RM has no significant variation to different parameter settings. This coincides with the discussion in Section 4.3: different parameter settings of IntSRM will generate almost the same quality after different numbers of interactive steps. We also find that the time cost of them has no much difference.

7 RELATED WORK

A host of works have been done on Schema Matching [31] and Record Matching [14]. While typical SM approaches are based on the similarity (or semantic correlation) between Attribute Names [10], [15], or Attribute Value Sets (i.e., instance-based SM) [2], [24] or combination of the two [12], typical RM methods measure the similarities between either key attribute values [4] or non-key attribute values [38]. Excellent surveys to the two problems can be found in [31] and [14], respectively.

For SM, our work is closely related to instance-based SM method, which has been recognized as an effective approach due to its robustness for matching heterogenous schemas [25], [27], [28]. Generally, the instance-based SM method leverages the classified instance data to process SM by measuring the similarity between sets of annotated instances (e.g., the construction of links between attributes based on the co-occurrence of instances). The basic idea of instancebased SM method is that the more significant overlap among common instances, the more relevant the two attributes are. The challenge here lies on how to define the significance for the overlap. [16] determines the similarity of two attributes by executing a pair-wise comparison of instance values using a similarity function, while [18] defines the similarity of two attributes by considering both the specificity and generalization of instances when two attributes are linked. In this paper, we propose a new SM rule by measuring the similarity of two attributes between the sets of linked record-pairs. More advanced, our method tends to link attributes with more explicit evidence (i.e., linked record-pair should have similar values under the same attribute), even there are just a few overlapped instances.

In addition, instance-based SM methods [18], [23] spend much time on comparing two set of attribute values. To save the time cost, some methods calculate the similarity between two attributes based on selected small subsets of attribute values that are generated from the two large original attribute value sets respectively. In [13], [23], hashing method and filtering strategies were proposed to improve the scalability of the instance-based schema matching. In this paper, we extend the q-gram index [9], [22], [37] to multiple pairs of attributes scenario, and split potential matched record-pairs between the two tables into (possibly overlapped) blocks such that matching-record-pairs are only identified within every block. At each SM and RM step, we only process those blocks that will bring the maximum benefits to the future interaction with the minimum cost.

Efficiency is a more serious problem for RM since there are usually many more records than attributes in databases. So far, various techniques have been proposed to reduce the overhead of RM, including Q-Grams together with inverted indices [34], prefix-based pruning techniques [35], batchbased matching techniques [7]. However, these techniques are only applied on the key attributes. A recent work [38] uses both key and non-key attributes for RM and relies on a special decision-tree to improve the efficiency of RM. In this paper, we do RM based on both key and non-key attribute values under the linked attributes, which outperforms previous methods on accuracy. Meanwhile, inspired by the previous methods, we extend the indices over all attributes to improve the efficiency of RM and the interaction process, which has been demonstrated to significantly reduce time cost in our experimental results.

There are also works on using structural/relational similarity of records for SM or RM. [3] and [33] resolve related entities collectively and combine attribute similarity with relational evidence to improve the quality of RM results. In [29], a graph based on attributes, attribute types and their relations is built to capture inherent structural information of the table, which can be utilized in SM. These approaches can improve the integration quality on data sets that have strong structural/relational information to utilize. However, the improvement of integration quality will be limited when data sets have weak structural information.

Recently, external domain knowledge and human interventions are also employed to improve the quality of SM [12] or RM [21]. These methods get external knowledge from either external knowledge base like wikipedia, or crowd workers, and use these external knowledge to label some attributes or instances and then extract effective features for better SM or RM.

So far, all existing efforts take SM and RM as independent steps in data integration. There has been little discussion about the interaction between them. A similar interesting study has been conducted on the interaction problem between RM and Data Cleaning [17], but the problem setting and the key challenges in that work are different from ours in various aspects.

8 **CONCLUSIONS AND FUTURE WORK**

In this paper, we study the interaction between SM and RM by performing them alternately in data integration. Extensive experiments on three data collections show that the combination and interaction between SM and RM significantly outperforms previous works that conduct SM and RM separately. Based on special-designed indices, we reduce around 90 percent overhead of the interaction algorithm.

Nonetheless, our approach has its own limitations: it can only work well with the case that one attribute (record) in a table matches with no more than one attribute (record) in the other table. As a future work, we will consider to extend our approach to deal with multiple matches. Our future work also includes addressing the problem when data is too large to be loaded into memory at a time.

ACKNOWLEDGMENTS

This research is partially supported by the Natural Science Foundation of China (Grant Nos. 61303019, 61402313, 61402312, 61472263, 61572336, 61572335, 61532018, 61502324), the King Abdullah University of Science and Technology, the Australian Research Council (Grants No. DP120102829), the Natural Science Foundation of Jiangsu (Grant No. BK20151223), the Postdoctoral scientific research funding of Jiangsu Province (No. 1501090B), and the National Postdoctoral Funding (Nos. 2015M581859, 2016T90493).

REFERENCES

- E. Agichtein and L. Gravano, "Snowball: Extracting relations from [1] large plain-text collections," in Proc. 5th ACM Conf. Digit. Libraries, 2000, pp. 85-94.
- P. A. Bernstein, J. Madhavan, and E. Rahm, "Generic schema [2] matching, ten years later," Proc. VLDB Endowment, vol. 4, no. 11, pp. 695–701, 2011. I. Bhattacharya and L. Getoor, "Collective entity resolution in rela-
- [3] tional data," ACM Trans. Knowl. Discovery Data, vol. 1, no. 1, 2007, Art. no. 5.
- [4] M. Bilenko and R. J. Mooney, "Adaptive duplicate detection using learnable string similarity measures," in Proc. 9th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 2003, pp. 39-48.
- C. M. Bishop, Pattern Recognition and Machine Learning. New York, [5] NY, USA: Springer, 2006.
- R. A. Brualdi, Introductory Combinatorics. Amsterdam, The Nether-[6] lands: North-Holland, 1992.
- A. Chandel, P. Nagesh, and S. Sarawagi, "Efficient batch top-k [7] search for dictionary-based entity recognition," in Proc. 22nd Int. Conf. Data Eng., 2006, pp. 28-28.

- P. Christen, "A survey of indexing techniques for scalable record [8] linkage and deduplication," IEEE Trans. Knowl. Data Eng., vol. 24, no. 9, pp. 1537–1555, Sep. 2012.
- P. Christen, "Towards parameter-free blocking for scalable record [9] linkage," Dept. Comput. Sci., Faculty Eng. Inf. Technol., Australian Nat. Univ., Canberra, Australia, Tech. Rep. TR-CS-07-03, 2007.
- [10] C. Comito, S. Patarin, and D. Talia, "A semantic overlay network for p2p schema-based data integration," in Proc. 11th IEEE Symp. Comput. Commun., 2006, pp. 88-94.
- [11] D. Dey, V. S. Mookerjee, and D. Liu, "Efficient techniques for online record linkage," IEEE Trans. Knowl. Data Eng., vol. 23, no. 3, pp. 373-387, Mar. 2011.
- [12] R. Dhamankar, Y. Lee, A. Doan, A. Halevy, and P. Domingos, "iMAP: Discovering complex semantic matches between database schemas," in Proc. ACM SIGMOD Int. Conf. Manage. Data, 2004, pp. 383–394.
- [13] S. Duan, A. Fokoue, O. Hassanzadeh, A. Kementsietsidis, K. Srinivas, and M. J. Ward, "Instance-based matching of large ontologies using locality-sensitive hashing," in Proc. 11th Int. Conf. Semantic Web-Part I, 2012, pp. 49-64.
- [14] A. K. Elmagarmid, P. G. Ipeirotis, and V. S. Verykios, "Duplicate record detection: A survey," IEEE Trans. Knowl. Data Eng., vol. 19, no. 1, pp. 1-16, Jan. 2007.
- [15] D. W. Embley, D. Jackman, and L. Xu, "Multifaceted exploitation of metadata for attribute match discovery in information integration," in Proc. Int. Workshop Inf. Integr. Web, 2001, pp. 110-117.
- [16] D. Engmann and S. Massmann, "Instance matching with COMA +," in Proc. BTW Workshops-Model Manage. Metadaten-Verwaltung, 2007, pp. 28-37.
- W. Fan, S. Ma, N. Tang, and W. Yu, "Interaction between record [17] matching and data repairing," J. Data Inf. Quality, vol. 4, no. 4, 2014, Art. no. 16.
- [18] A. Ferraram, A. Nikolov, and F. Scharffe, "Data linking for the Semantic Web," in Semantic Web: Ontology and Knowledge Base Enabled Tools, Services, and Applications. Hershey, PA, USA: Idea Group Inc., 2013, p. 169.
- [19] A. Gal, "Managing uncertainty in schema matching with top-k schema mappings," in *Journal on Data Semantics VI*. Berlin, Germany: Springer-Verlag, 2006, pp. 90–114.
- [20] F. Giunchiglia, M. Yatskevich, and P. Shvaiko, "Semantic matching: Algorithms and implementation," in Journal on Data Semantics IX. Berlin, Germany: Springer, 2007, pp. 1–38.
- [21] C. Gokhale et al., "Corleone: Hands-off crowdsourcing for entity matching," in Proc. ACM SIGMOD Int. Conf. Manage. Data, 2014, pp. 601–612.
- [22] L. Gravano, P. G. Ipeirotis, H. V. Jagadish, N. Koudas, S. Muthukrishnan, and D. Srivastava, "Approximate string joins in a database (almost) for free," in Proc. 27th Int. Conf. Very Large Data Bases, 2001, vol. 1, pp. 491–500. [23] O. Hassanzadeh, et al., "Discovering linkage points over web
- data," Proc. VLDB Endowment, vol. 6, no. 6, pp. 445-456, 2013.
- [24] T. Kirsten, A. Thor, and E. Rahm, "Instance-based matching of large life science ontologies," in Data Integration in the Life Sciences. Berlin, Germany: Springer, 2007, pp. 172-187.
- [25] J. Li, et al., "Diversity-aware retrieval of medical records," Comput. Ind., vol. 69, pp. 81-91, 2015.
- [26] W. Lin, R. Yangarber, and R. Grishman, "Bootstrapped learning of semantic classes from positive and negative examples," in Proc. 20th Int. Conf. Mach. Learn. Workshop Continuum Labeled Unlabeled Data, vol. 1, no. 4, p. 21, 2003.
- [27] R. Mao, H. Xu, W. Wu, J. Li, Y. Li, and M. Lu, "Overcoming the challenge of variety: Big data abstraction, the next evolution of data management for AAL communication systems," IEEE Commun. Mag., vol. 53, no. 1, pp. 42-47, Jan. 2015.
- [28] R. Mao, P. Zhang, X. Li, X. Liu, and M. Lu, "Pivot selection for metric-space indexing," Int. J. Mach. Learn. Cybern., vol. 7, no. 2, pp. 311–323, 2016.
- [29] S. Melnik, H. Garcia-Molina, and E. Rahm, "Similarity flooding: A versatile graph matching algorithm and its application to schema matching," in *Proc. 18th Int. Conf. Data Eng.*, 2002, pp. 117–128. [30] D. Pollard, "Quantization and the method of k-means," *IEEE*
- *Trans. Inf. Theory*, vol. 28, no. 2, pp. 199–205, Mar. 1982. [31] E. Rahm and P. A. Bernstein, "A survey of approaches to auto-
- matic schema matching," VLDB J., vol. 10, no. 4, pp. 334–350, 2001.
- [32] P. Shvaiko and J. Euzenat, "A survey of schema-based matching approaches," in Journal on Data Semantics IV. Berlin, Germany: Springer, 2005, pp. 146–171.

- [33] J. Tang, A. C. Fong, B. Wang, and J. Zhang, "A unified probabilistic framework for name disambiguation in digital library," *IEEE Trans. Knowl. Data Eng.*, vol. 24, no. 6, pp. 975–987, Jun. 2012.
- [34] E. Ukkonen, "Approximate string-matching with q-grams and maximal matches," *Theoretical Comput. Sci.*, vol. 92, no. 1, pp. 191– 211, 1992.
- [35] W. Wang, C. Xiao, X. Lin, and C. Zhang, "Efficient approximate entity extraction with edit distance constraints," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, 2009, pp. 759–770.
- [36] L. Wasserman, All of Statistics. New York, NY, USA: Springer, 2011.
- [37] C. Xiao, W. Wang, and X. Lin, "Ed-join: An efficient algorithm for similarity joins with edit distance constraints," *Proc. VLDB Endowment*, vol. 1, no. 1, pp. 933–944, 2008.
- *ment*, vol. 1, no. 1, pp. 933–944, 2008.
 [38] Q. Yang, et al., "NokeaRM: Employing non-key attributes in record matching," in *Proc. 16th Int. Conf. Web-Age Inf. Manage.*, 2015, pp. 438–442.



Binbin Gu is working toward the master's degree in the Research Center on Advanced Data Analytics, Soochow University, China. He visited Prof. Xiangliang Zhang at the King Abdullah University of Science and Technology for half a year in 2016. His research interests include knowledge fusion, information extraction, and NLP. He has published several papers at DASFAA. He is the external reviewer of several international conferences such as ADC and WAIM.



Zhixu Li received the BS and MS degrees from the Renmin University of China, in 2006 and 2009, respectively, and the PhD degree from the University of Queensland, in 2013. He is an associate professor in the Department of Computer Science & Technology, Soochow University, China. He used to work as a research fellow with the King Abdullah University of Science and Technology. His research interests include data cleaning, big data applications, information extraction, and retrieval.



Xiangliang Zhang received the PhD degree in computer science from INRIA-University Paris-Sud 11, France, in July 2010. She is an assistant professor and directs the Machine Intelligence and kNowledge Engineering Laboratory, King Abdullah University of Science and Technology. Her main research interests and experiences include diverse areas of machine learning and data mining.



An Liu received the PhD degree from both the City University of Hong Kong (CityU) and University of Science & Technology of China (USTC), in 2009. He is an associate professor in the Department of Computer Science & Technology, Soochow University. Prior to that in 2014, he was a senior research associate in the Joint Research Center, CityU and USTC. His research interests include security, privacy, and trust in emerging applications, cloud computing, and services computing.



Guanfeng Liu received the PhD degree in computer science from Macquarie University, Australia, in 2013. He is an associate professor in the School of Computer Science and Technology, Soochow University, China. His research interests include social network mining and trust. He has published more than 40 papers in the most prestigious journals and conferences. He was the PC chair of BDMS 2015.



Kai Zheng received the PhD degree in computer science from The University of Queensland, in 2012. He is a professor in the School of Computer Science and Technology, Soochow University. His research focus is to find effective and efficient solutions for managing, integrating and analyzing big data for business, and scientific and personal applications. He has been working in the area of spatial-temporal databases, uncertain databases, trajectory computing, socialmedia analysis, and bioinformatics. He has pub-

lished more than 60 papers in the highly referred journals and conferences such as SIGMOD, ICDE, EDBT, *The VLDB Journal*, ACM Transactions, and IEEE Transactions.



Lei Zhao received the PhD degree in computer science from Soochow University, in 2006. He is a professor in the School of Computer Science and Technology, Soochow University. His research focuses on graph databases, social media analysis, query outsourcing, parallel, and distributed computing. His recent research is to analyze large graph database in an effective, efficient, and secure way. He has published more than 100 papers including more than 20 published in well-known journals and conferences

such as ICDE, DASFAA, WISE, and the *Journal of Computer Science* and *Technology*.



Xiaofang Zhou is a professor of computer science with The University of Queensland. He is the head of the Data and Knowledge Engineering Research Division. He is a specially appointed adjunct professor under the Chinese National Qianren Scheme hosted by the Renmin University of China (2010-2013), and by Soochow University since July 2013, where he leads the Research Center on Advanced Data Analytics. He has been working in the area of spatial and multimedia databases, data quality, high perfor-

mance query processing, Web information systems, and bioinformatics, co-authored more than 250 research papers with many published in top journals and conferences. He is a senior member of the IEEE.

▷ For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.