

Query-Driven Sampling for Collective Entity Resolution

October 2, 2018

Abstract

Probabilistic databases play a preeminent role in the processing and management of uncertain data. Recently, many database research efforts have integrated probabilistic models into databases to support tasks such as information extraction and labeling. Many of these efforts are based on batch oriented inference which inhibits a realtime workflow. One important task is entity resolution (ER). ER is the process of determining records (mentions) in a database that correspond to the same real-world entity. Traditional pairwise ER methods can lead to inconsistencies and low accuracy due to localized decisions. Leading ER systems solve this problem by collectively resolving all records using a probabilistic graphical model and Markov chain Monte Carlo (MCMC) inference. However, for large datasets this is an extremely expensive process. One key observation is that, such exhaustive ER process incurs a huge up-front cost, which is wasteful in practice because most users are interested in only a small subset of entities.

In this chapter, we advocate pay-as-you-go entity resolution by developing a number of query-driven collective ER techniques. We introduce two classes of SQL queries that involve ER operators — selection-driven ER and join-driven ER. We implement novel variations of the MCMC Metropolis Hastings algorithm to generate biased samples and selectivity-based scheduling algorithms to support the two classes of ER queries. Finally, we show that query-driven ER algorithms can converge and return results within minutes over a database populated with the extraction from a newswire dataset containing 71 million mentions.

1 Query-Driven Entity Resolution Introduction

Entity resolution (ER) is the process of identifying and linking/grouping different manifestations (e.g., mentions, noun phrases, named entities) of the same real world object. It is a crucial task for many applications including knowledge base construction, information extraction, and question answering. For decades, ER has been studied in both database and natural language processing communities to link database records or to perform entity resolution over extracted mentions (noun phrases) in text.

ER is a notoriously difficult and expensive task. Traditionally, entities are resolved using strict pairwise similarity, which usually leads to inconsistencies and low accuracy due to localized, myopic decisions [39]. More recently, collective entity resolution methods have achieved state-of-the-art accuracy because they leverage relational information in the data to determine resolution jointly rather than independently [5]. However, it is expensive to run collective ER based on probabilistic graphical models (GMs), especially for cross-document entity resolution, where ER must be performed over millions of mentions.

In many previous approaches, collective ER is performed exhaustively over all the mentions in a data set, returning all entities. Researchers have developed new methods to perform large-scale cross-document entity resolution over parallel frameworks [33, 39]. However, in many ER applications, users are only interested in one or a small subset of entities. This key observation motivates query-driven ER, an alternative approach to solving the scalability problem for ER.

Compared to previous ER models and algorithms, query-driven techniques in this chapter scale to data sets that are in many cases three orders of magnitude larger. Moreover, the ER model in this chapter is general enough to take both bibliographic records and mentions extracted from unstructured text. Query-driven ER techniques over GMs can also be generalized for other applications to perform query-driven inference.

This work follows a line of research on implementing ML models inside of databases [18, 22, 38]. Researchers use factor graphs because this flexible representation works well with other machine learning algorithms. ER is ubiquitous and an important part of many analytic pipelines; a probabilistic database implementation is natural.

In this chapter, we first introduce SQL-like queries that involve ER operations. These ER operators are an SQL comparison operator (i.e., ER-based equality) that returns true if two mentions map to the same entity.

Factor Graphs, a type of GM, are used to model the collective entity resolution over extracted mentions from text. Using this ER based comparison operator, users can pose selection queries to find all mentions that map to a single entity or pose join queries to find mentions that map to the subset of entities that they are interested in resolving.

Because exhaustive ER is expensive it is common to use blocking techniques to partition the data set into approximately similar groups called canopies. Query-driven ER in this chapter differs from blocking in two important ways: 1) deterministic blocks are replaced by a pairwise distance-based metric, and 2) blocks (or canopies) are implicit to the query-driven ER data set and do not have to be created in advanced. The latter point, implicit blocking, is realized using a data structure created based on the similarity to a query mention. This data structure allows parameters to include or remove mentions from the working data set. This property is similar to the iterative blocking technique [37], which is shown to improve ER accuracy. Such an approach can dramatically amortize the overall ER cost suitable for the pay-as-you-go paradigm in dataspace [25].

To support ER driven by queries, we develop three sampling algorithms for MCMC inference over graphical models. More specifically, instead of a uniform sampling distribution, we sample on a distribution that is biased to the query. We develop a query-driven sampling techniques that maximizes the resolution of the target query entity (target-fixed) and biases the samples based on the pairwise similarity metric between mentions and query nodes (query-proportional). We also introduce a hybrid method that performs query-proportional sampling over a fixed target. We develop two optimizations to the query-proportional and hybrid methods to model the similarity and dissimilarity between the mentions and the query entity, i.e., attract and repel scores. In the first target-fixed algorithm, we adapt the samples to resolve the query entity. The second query-proportional algorithm, selects mentions based on their probabilistic similarity to the query entity. The third hybrid algorithm combines the two approaches. A summary of approaches can be found in Table 3.

When a user is interested in resolving more than one entity we employ multi-node ER techniques. To implement multi-node ER queries, single-node ER techniques may be naively performed iteratively to resolve one entity at a time. However, such an algorithm can lead to un-optimized resource allocation if the same number of samples is generated for each target entity, or low throughput if one of the entities has a disproportionately low convergence rate. To alleviate this problem, we present three multi-query ER algorithms that schedule the sample generation among query nodes in order to improve overall convergence rate.

In summary, the contributions of this chapter are the following:

- We define a query-driven ER problem for cross-document, collective ER over text extracted from unstructured data sets;
- We develop three single-node algorithms that perform focused sampling and reduce convergence time by orders-of-magnitude compared to a non-query-driven baseline (Section 4). We develop two influence functions that use attract and repel techniques to grow or shrink query entities (Section 5.1);
- We develop scheduling algorithms to optimize the overall convergence rate of the multi-query ER (Section 5.2). The best scheduling algorithm is based on selectivity of different target entities (Section 5.3).

The results show that query-driven ER algorithms is a promising method of enabling realtime, ad-hoc, ER-based queries over large data sets. Single node queries of different selectivity converge to a high-quality entity within 1-2 minutes over a newswire data set containing 71 million mentions. Experiments also show that such real-time ER query answering allows users to iteratively refine ER queries by adding context to achieve better accuracy (Section 6).

2 Query-Driven Entity Resolution Preliminaries

In this section we present a foundation of concepts discussed in this chapter. We start with an introduction of factor graphs then discuss sampling techniques over this model. Finally, we formally introduce state-of-the-art entity resolution approaches and explain the origin.

2.1 Factor Graphs Graphical models are a formalism for specifying complex probability distributions over many interdependent random variables. Factor graphs are bipartite graphical models that can capture arbitrary relationships between random variables through the use of factors [21]. As depicted in Figure 1, links always connect random variables (represented as circles) and factor nodes (represented as black squares). Factors are

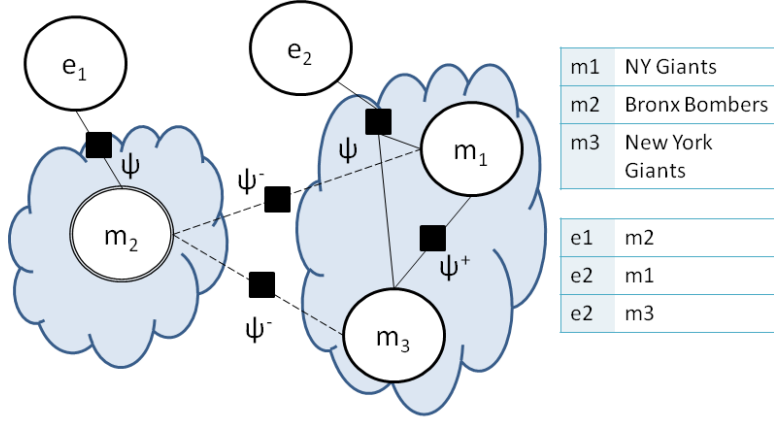


Figure 1: Three node factor graph. Circles (random variables) with m_i represent mentions and those with e_i represent entities. Clouds are added for visual emphasis of entity clusters

functions that take as input the current setting of connected random variables, and output a positive real-valued scalar indicating the compatibility of the random variables settings. The probability of a setting to all the random variables is a normalized product of all the factors. Intuitively, the highest probability settings have variable assignments that yield the highest factor scores.

We use factor graphs to represent complex entity resolution relationships. Nodes (random variables) may correspond to mentions of people, places and organizations in documents. Nodes also represent the random variables that correspond to groups of mentions (entities), these nodes are accompanied by clouds in Figure 1. The factors between mentions and entities give us a sound representation for many possible states. The factor graph model also gives us a simple mathematical expression of the relationship.

Formally, a factor graph $\mathcal{G} = \langle \mathbf{x}, \psi \rangle$ contains a set of random variables $\mathbf{x} = \{x_i\}_1^n$ and factors $\psi = \{\psi_i\}_1^m$. Each factor ψ_i maps the subset of variables it is associated with to a non-negative compatibility value. The probability of a setting ω among the set of all possible settings Ω occurring in the factor graph is given by a probability measure:

$$\pi(\omega) = \frac{1}{Z} \sum_{x \in \omega} \prod_{i=1}^m \psi_i(x^i), \quad Z = \sum_{\omega \in \Omega} \sum_{x \in \omega} \prod_{i=1}^m \psi_i(x^i)$$

where x^i is the set of random variables that neighbor the factor $\psi_i(\cdot)$ and Z is the normalizing constant.

Querying graphical models produces the most likely setting for the random variables. A query on a factor graph is defined as a triple $\langle x_q, x_l, x_e \rangle$ where x_q is the set of nodes in question, x_l is a set of latent nodes (entities) that are marginalized and x_e is a set of evidence nodes (observed mentions). A query task is a sum over the all latent variables and the maximization of the query probability. A query over the factor graph is defined as

$$\mathcal{Q}(x_q, x_l, x_e, \pi) = \operatorname{argmax}_{x_q} \sum_{v_l \in x_l} \pi(x_q \cup v_l \cup x_e).$$

To obtain the best setting of the queries in question, inference is required.

Several methods exist for performing inference over factor graphs. The entity resolution factor graph, being pairwise, is dense and highly connected. This property suggests the best methods for inference are Markov Chain Monte Carlo (MCMC) methods; in particular, we use a Metropolis Hastings variant [21]. We refer the reader to our previous work for a detailed discussion on inference over factor graphs and a deviation of the technique [40].

2.2 Inference over Factor Graphs Several methods exist for performing inference over factor graphs. The entity resolution factor graph, being pairwise, is dense and highly connected. This property suggests the best methods for inference are Markov Chain Monte Carlo (MCMC) methods; in particular, we use a Metropolis Hastings variant [21].

The idea of MCMC-MH is to propose modifications to a current setting and use the model to decide whether to accept or reject the proposed setting as a replacement for the current settings. When the models are being scored only the factors touching nodes with changed values, the Markov blanket, needs to be recomputed. We accept or reject changes so the model can iteratively proceed to an optimal setting.

More formally, consider an MCMC transition function $T : \Omega \times \Omega \rightarrow [0, 1]$ where given the current setting ω we can sample a subsequent setting ω' .

The probability of accepting a transition given a graphical model distribution π is:

$$(2.1) \quad A(\omega, \omega') = \min \left(1, \frac{\pi(\omega')T(\omega, \omega')}{\pi(\omega)T(\omega', \omega)} \right).$$

Additionally, the intractable partition function Z is canceled out, making sample generation inexpensive. This property allows us to calculate the probability of accepting the next state by simply computing the difference in score between the next and current state [40].

We say the algorithm converges when a steady state is reached.¹ Intelligently sampling next states decreases the time to convergence. Convergence in MCMC is difficult to verify [10], we discuss convergence estimation in Section 6.1.

2.3 Cross-Document Entity Resolution Cross-document ER is the problem of clustering mentions that appear across independent sets of documents into groups of mentions that correspond to the same real world entity. These ER tasks typically assume a set of preprocessed documents and perform linking across documents [4, 33]. The scale of the cross-document ER problem is typically several orders of magnitude more than intra-document ER. There are no document boundaries to limit inference scope and all entity mentions may be distributed arbitrarily across millions of documents.

To model cross-document ER, let $\mathcal{M} = \{m_1, \dots, m_{|\mathcal{M}|}\}$ be the set of mentions in a data set. Each mention m_i contains a set of attribute-value data points. Let $\mathcal{E} = \{e_1, \dots, e_{|\mathcal{E}|}\}$ represent the set of entities where each e_i contain zero or more mentions. Note, we assume the maximum number of entities is no more than the number of mentions and no less than 1. Each mention may correspond to a unique entity or all mentions may correspond to a single entity.

The baseline method of entity resolution is a straight-forward application of the MCMC-MH algorithm. We show pseudo code for the baseline method in Algorithm 1.

Algorithm 1 The baseline entity resolution algorithm using Metropolis-Hastings sampling

INPUT: A set of unresolved entities \mathcal{E} each with one mention m .

INPUT: A positive integer *samples*.

OUTPUT: A set of resolved entities \mathcal{E} .

```

1: while samples-- > 0 do
2:    $e_i \sim_u \mathcal{E}$ 
3:    $e_j \sim_u \mathcal{E}$ 
4:    $m \sim_u e_i$ 
5:    $\mathcal{E}' \leftarrow \text{MOVE}(\mathcal{E}, m, e_j)$ 
6:   if  $\text{SCORE}(\mathcal{E}) < \text{SCORE}(\mathcal{E}')$  then
7:      $\mathcal{E} \leftarrow \mathcal{E}'$ 
8:   end if
9: end while
   return  $\mathcal{E}$ 

```

Algorithm 1 takes as input a set of entities \mathcal{E} and *samples* which is the number of iterations of the algorithm or a function to estimate convergence. The algorithm samples two entities from the entity set and moves one

¹We refer to literature for a more detailed description of convergence [40].

random² mention into the other entity. After the move, the algorithm checks for an improvement in the overall score of the model. If the model score improves, the changes are kept, otherwise the proposed changes are ignored. The SCORE function sums the weights of all the edges in the given entity to obtain a value for the model. This is equivalent to the probability of the setting $\pi(\cdot)$ as described in Section 2.1.

Blocking. Blocking or canopy generation is a preprocessing technique to partition large amounts of data into smaller chunks, or blocks of items that are likely to be matches [26]. Blocking can use simple and fast techniques such as sorting based on attributes or more advanced techniques that map similar items onto a vector space [11, 32, 37].

In this chapter, we use two methods of blocking. First, we use an approximate string match over all the mentions in the database. To perform the approximate string filter we use a q -grams technique over all the mentions in the database. This method creates an inverted index for each mention in the database so a query can be performed to look for all words that contain a sufficient number of matching q -grams. This gives us a fast high-recall filter over many records [17].

The second is an implicit blocking structure created by computing the influence a query node has on the other nodes in the data set (see Section 5.1). This method uses an estimate of the distance between the query nodes and the candidate mentions to prioritize samples.

3 Query-Driven Entity Resolution Problem Statement

In this section, we formally define the problem of query-driven ER. We use an SQL-like formalism to model traditional and query-driven entity resolution.

In a probabilistic database, let a Mentions table contain all the extracted mentions from a text corpus. Its column `entityp` represents the probabilistic latent entity labels; they contain a mapping but that mapping may not represent the current state. The People table holds a watchlist of mentions and relevant contextual information. The context column is an abstract place holder for text data or richer schemas. This model only assumes there is a master column, the realization of the context column is flexible and implementation dependent.

```
Mentions(docID, startpos, mention, entityp, context)
People(peopleID, mention, entityp, context)
```

We also define a user-defined function `coref_map` that performs maximum a posteriori (MAP) inference on the latent entity^p random variables. The function takes two instances of mentions with at least one being from a probabilistic table such as the Mention table. When the query is executed the `coref_map` function returns true if the mentions referenced are coreferent. Following, we describe the traditional exhaustive ER task as well as the single- and multi-node query-driven ER queries.

Exhaustive The goals of traditional entity resolution is to cluster all mentions in a data set. All the mentions clustered inside each entity are coreferent with each other and not entity with any mention that is a part of a different entity cluster. The process of exhaustive ER can be modeled as a self-join database query where each mention is grouped into coreferent clusters. In Algorithm 2 we create a view displaying the results of a resolved query.

Algorithm 2 Example exhaustive entity resolution query that creates a database view

```
CREATE VIEW CorefView AS
SELECT m.docID, m.startpos, m.mention, m2.mention
FROM Mentions m, Mentions m2
WHERE coref_map(m.*, m2.entityp), m2.mention, m2.context)
```

To obtain unique entity clusters, we can perform an aggregation query over the CorefView. In Figure 3 we see an example of the result of traditional entity resolution.

²Given a set X , the function $x \sim_u X$ makes a uniform sample from the set X into a variable x .

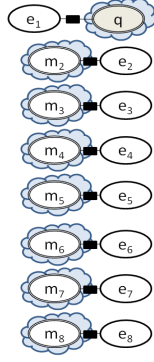


Figure 2: A possible initialization for entity resolution

Single-node Query In the ER task, we may only be interested in the mentions of one entity. We represent this entity with a template mention, or as a query node q . Single-node entity resolution is modeled as a selection query with a where-clause that includes the template mention q and returns only the mentions that are members of the entity cluster that contains the sample mention. Given a template mention q and its context $q.context$, Algorithm 3 we show the single-node query based on an example in Section 4.2.

Algorithm 3 Single query-node driven entity resolution query

```
SELECT m.docID, m.startpos, m.mention
FROM Mentions m
WHERE coref_map(m.*, m.entity^p), q, q.context)
```

Here we add parameters to the `coref_map` function that contain the specific query and its context. It performs ER over the mentions table but only returns an affirmative value if the labels for the entity cluster match the query node. For example, if the template mention q was ‘Mark Zuckerberg’, and the query context were keywords such as ‘facebook’ and ‘ceo’, the only returned mentions will be those that represent Mark Zuckerberg the facebook founder. This is similar to a ‘facebook’ approximate string search. The emphasis of this chapter is optimizing this function so while performing ER we perform less work compared an exhaustive query.

Multi-Query In many cases, a user may be interested in a watchlist of entities. Watchlist is a subset of the larger Mention set. This is common for companies looking for mentions of its products in a data set. In this case, mention are only clustered with the entities represented in the watchlist. Algorithm 4 is an example of a join-query between the Mentions table and the People table.

Algorithm 4 Multi-query between the Peoples watch list tablse and the full mentions set

```
SELECT m.docID, m.startpos, m.mention, q
FROM Mentions m, People q
WHERE coref_map(m.*, m.entity^p), q, q.context)
```

This function combines a watch list of terms and performs ER with respect to the specific examples in the watchlist. The multi-query method uses scheduling to perform inference, or a fuzzy equal, over each mentions. In Section 4.3 we propose scheduling algorithms so multi-query node ER gracefully manage multiquery workloads.

4 Query-Driven Entity Resolution Algorithms

Query-driven ER is an understudied problem; in this section we describe our approach to query-driven ER with one entity (single-query ER) and with multiple entities (multi-query ER). First, we give a graphical intuition of query-driven ER algorithms.

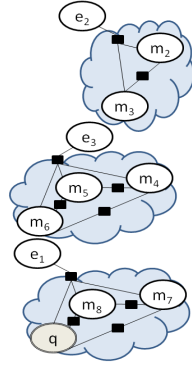


Figure 3: The correct entity resolution for all mentions

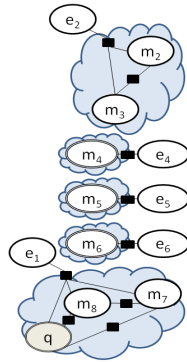


Figure 4: The entity containing q is internally coreferent; the other entities are not correctly resolved

4.1 Intuition of Query-Driven ER In this section, we remind the reader of the query-driven ER task with a formal definition. Each ER task is given a corpus \mathcal{G} and a set of entity mentions $\mathcal{M} = \{m_1, \dots, m_{|\mathcal{M}|}\}$ extracted from the \mathcal{G} . A user may supply set of query nodes $Q = \{q_1, \dots, q_{|Q|}\}$. Each q_i , also called a query template, may be a member of \mathcal{M} or a manually declared mention that is appended to the set of mentions. For each node $q_i \in Q$, the task of ER is to compute the set of mentions $E = \{e_1, \dots, e_{|Q|}\}$ that only contain mentions that are coreferent with the query node,

$$e_{q_i} = \{m_i | m_i \in \mathcal{M}, \text{QDER}(\mathcal{M}, m_i, q_i)\}.$$

In Section 4.2, we describe implementations of the QDER algorithm for $|Q| = 1$. In Section 4.3, we describe techniques of scheduling the ER task for the general case of $|Q| > 1$.

Fundamentally, the ER algorithm generates a graphical model and makes new state proposals (jumps) to reach the best state (see Section 2). The query-driven algorithms in this section use a query node to facilitate more sophisticated jumps. By making smart proposals we expect faster convergence to an accurate state. As a note to the reader, a summary of query-driven algorithms can be found in Table 3.

Figure 4 shows an initial configuration and acceptable query-driven entity resolution solutions. An example initial state of this algorithm is shown in Figure 2 — each mention is initially assigned to separate entities. Alternatively, the model may be initialized randomly, or in an arrangement from a previous entity resolution output or with all mentions in one entity. Figure 3 is the full resolution for the data set; each mention is correctly assigned to its entity cluster. Figure 4 is a result that was resolved with query-driven methods and is a partially resolved data. Because the entity containing the query node is completely resolved the solution is acceptable.

4.2 Single-Node ER Single-node ER algorithms are the class of algorithms that resolve a single query-node as discussed in Section 3. In particular, the target-fixed ER algorithm aims to focus a majority of the proposals on resolving the query entity. The algorithm fixes the query node as the target entity and then randomly selecting a source node to merge into the entity of the target query node. This focus on building the query entity in this type of importance sampling means the query entity should be resolved faster than if we sampling each entity uniformly.

A query-driven ER algorithm that only selects the query-node as the target entity during sampling will create errors because such an algorithm is unable to remove erroneous mentions from the query entity. To prevent these errors, we allow the algorithm to occasionally back out of poor decisions, that is, it makes non-query specific samples. Shown in Algorithm 5, target-fixed entity resolution adapts Algorithm 1 but it allows parameters to specify the proportion of time the different sampling methods are selected.

In addition to the input mentions \mathcal{E} from Algorithm 1, target-fixed entity resolution takes as input a query node q . The output of the algorithm is a resolved query entity and other partially resolved entities.

For each sampling iteration the algorithm can make two decisions. The sampler may propose to merge a random source node that is not already a member of the query entity into the target query entity. Alternatively, the algorithm merges a random node with a random entity.

On lines 3 to 6 the algorithm takes a uniform sample from the list of entities. If the sampled entity is the same as the query entity it tries again and samples a distinct entity. A node is drawn from this entity. The probability of this block being entered is τ_α . Lines 7 to 10 are entered with a probability $(1 - \tau_\alpha)$. This block performs a random entity assignment in the same manner as Algorithm 1. This block offsets the aggressive nature of the target-fixed algorithm by probabilistically backing out of any bad merges. Finally, the block starting from line 12 to line 15 scores the new arrangement and accepts if this improves the model score. We discuss parameter settings in Section 5.5.

Example Take the synthetic mention set \mathcal{M} shown in Table 1 and a query node q , the baseball team ‘New York Yankees’, in Table 2. This is the result of the approximate match of query q over a larger data set (blocking). The mentions of \mathcal{M} may be initialized by assigning each mention to its own entity. After a successful run of traditional entity resolution the set of entities clusters are

$$\{\langle q, m_2, m_4, m_6 \rangle, \langle m_1, m_3 \rangle, \langle m_5 \rangle\}.$$

For query-driven scenario the only entity we are interested in is $\langle q, m_2, m_4, m_6 \rangle$. Each mention in this query entity is an alias for the ‘New York Yankees’ baseball team. The other two mentions represent the ‘New York Giants’ football team and the ‘Brooklyn Dodgers’ baseball team respectively.

Algorithm 5 Target-fixed entity resolution algorithm

Input: A query node q .
 A set of entities \mathcal{E} each with one mention m .
 A positive integer *samples*.

Output: A set of resolved entities \mathcal{E}' .i

```
1:  $\mathcal{E}' \leftarrow \mathcal{E} \cup q$ 
2: while samples-- > 0 do
3:    if RANDOM() <  $\tau_\alpha$  then
4:      $e_i \sim_u \mathcal{E}'$ 
5:      $e_j \leftarrow q.entity$ 
6:      $m \sim_u e_i$ 
7:    else
8:      $e_j \leftarrow \{e | \exists e, e \in \mathcal{E}', e \neq q.entity\}$ 
9:      $e_i \leftarrow \{e | \exists e, e \in \mathcal{E}', e \neq e_j\}$ 
10:     $m \sim_u e_i$ 
11:    end if
12:     $\mathcal{E}'' \leftarrow \text{MOVE}(\mathcal{E}', m, e_j)$ 
13:    if SCORE( $\mathcal{E}'$ ) < SCORE( $\mathcal{E}''$ ) then
14:      $\mathcal{E}' \leftarrow \mathcal{E}''$ 
15:    end if
16: end while
    return  $\mathcal{E}'$ 
```

Table 1: Mentions sets \mathcal{M} from a corpus

id	Mention	...
m_1	NY Giants	...
m_2	Bronx Bombers	...
m_3	New York Giants	...
m_4	Yankees	...
m_5	Brooklyn Dodgers	...
m_6	The Yanks	...

Table 2: Example query node q

id	Mention	...
q	New York Yankees	...

The target-fixed algorithm attempts to merge nodes with the query entity one mention at a time and the merge is accepted if it improves the score of the overall model. We can see in the example that a merge of m_1 and m_3 may improve the overall model because they have similar keywords but one refers to the query entity and the other to different football team. The target-fixed algorithm can correct this type of error by probabilistically backing out of errors by moving mentions in the query node to a new entity as show in line 7 to line 10 of Algorithm 5.

4.3 Multi-query ER A user may want to resolve more than one query entity, that is, she may be interested in resolving a watch list of entities over the data set. To support multiple queries, first merge the canopies of each query node in the watch list to obtain a subset of the full graphical model containing only the nodes similar to query nodes. To resolve the entities we can use query-proportional methods iteratively over each query node. We define two classes of schedules, namely, static and dynamic.

Static schedules are formulated before sampling while dynamic schedules are updated in response to estimated convergence. The two static schedules we develop are random and selectivity-based. In random scheduling each query node from the watch list is selected in a round robin style. Selectivity-based scheduling is a method of ordering multi-query samples to schedule proposals in proportion to the selectivity of the query node. Selectivity, in this case, is defined as the number of mentions retrieved using an approximate match of the data set, or the query node’s contribution to the total new graphical model. For example, the selectivity of our query node q in Table 2 the selectivity is simply the size of \mathcal{M} , shown in Table 1.

Random-based scheduling method performs well if all query nodes come from similar selectivity. Otherwise, if the selectivity of each query node vary, one query node may require more sampling compared to the others. If one query node needs a lot of samples to converge, it may take the whole process a long time to complete and cycles may be wasted on other nodes that have already converged.

In addition to scheduling samples in proportion to their selectivity, we can schedule samples dynamically, depending on the progress of each query entity. To perform dynamic scheduling we need to know how each query entity is progressing towards convergence. To estimate the running convergence we do not use standard techniques in literature because scheduling needs to occur before the model is close to convergence [10]. Instead, we estimate the convergence by measuring the fraction of accepted samples over the last N samples of each query in the watch list. The two dynamic scheduling algorithms are closest-first and sampling the farthest-first. In closest-first we queue up the query node that has the lowest positive average number of accepted nodes over the last N proposals. This scheduling method performs inference for the node that is closest to being resolved so it can move on to other nodes. Alternatively, the farthest-first algorithm schedules the node that has the highest convergence rate. This scheduling algorithm makes each query entity progress evenly.

5 Optimization of Query-Driven ER

The previous ER techniques aggressively attempt to resolve the query entity. However, if the query node is not representative of the query items performance of target-fixed ER can lead to undesirable results. We do not explore this trade-off; we assume users can select representative query nodes. In this section, we introduce optimizations to create approximate query-driven samples based on the query node. We first discuss the influence function that is used to make query-driven proposals. We then discuss the attract and repel versions of the influence function followed by two new algorithms. We end with implementation details and a summary of our query-driven algorithms.

5.1 Influence Function: Attract and Repel To retrieve nodes from a graphical model that is similar to a query node we employ the notion of influence. Our assumption is that nodes that are similar have a high probability of being coreferent. An influence trail score between two nodes in a graphical model can be computed as the product of factors along their active trail as defined in literature [40]. For a node $m_i \in \mathcal{M}$ and the query node $q \in \mathcal{M}$ the influence of m_i on the query node is defined as:

$$\mathcal{I}(m_i, q) = \sum_{j \in \mathcal{F}} w_j \psi_j(m_i, q)$$

where \mathcal{F} is the world of pairwise features and the feature weight and log-linear function are, respectively, w_j and ψ_j . The influence function \mathcal{I} is an implementation of this trail score.

The influence function takes a set of entities — or the equivalent GM — and a query node q as parameters. The parameters to an influence function can be over the whole database or a canopy. Over several invocations of the function, \mathcal{I} returns mentions from the graphical model with a frequency proportionate to their influence on q . If a mention has little or no influence, the influence acts as a blocking function, infrequently returning the mention. Recall influence is the distance active trail distance to query node. To implement the influence function we build a data structure based on an algorithm by Vose [36], hereafter referred to as a Vose structure.

The input mentions to the blocking algorithms may result in high or low quality canopies. A high quality canopy means most of the mentions in the canopy are associated with the query node. Low quality canopies, which are more common, corresponds to only a small number of mentions being associated with the query node. When initializing query-driven algorithms the canopy quality is important for determining what algorithm to use.

The attract method initializes each mention in the canopy in its own entity, and then mentions are merged until the convergence. The target-fixed algorithm discussed in Section 4.2 is explained using this method. The attract method works well for low quality canopies, or canopies that require a small number of items to merge. Conversely, the repel method works well with high quality canopies or when most items in a canopy belong to the query entity.

The repel method initializes each mention in the canopy into a single entity. Then proposals are made to remove mentions from the entity so we are left with only the nodes in the query entity. We discuss this method using the hybrid algorithm in Section 5.3. To build an influence function for the repel method we can use the same method and we only need to normalize and invert the influence scores. We refer to this as co-influence or $\tilde{\mathcal{I}}$.

5.2 Query-proportional ER In the query-proportional sampling algorithm, on every iteration, the source mention and target entity are selected in proportion to its distance to the query entity. Instead of focusing solely on the query entity, this algorithm prioritizes samples using a measure that represents probability of a mention being coreferent with the query entity.

That is, each node p in the graphical model \mathcal{G} is selected on the active trail between itself and the query node q . This algorithm merges nodes that are similar to the query node with an increased frequency.

Before query-proportional sampling, a data structure for \mathcal{I} is created. The \mathcal{I} influence structure takes a query node q and the global graphical model \mathcal{E} then returns a sampled mention. As \mathcal{I} is called multiple times, the distribution of the nodes returned is proportional to their influence. Algorithm 6 describes the query-proportional algorithm.

Algorithm 6 Query-proportional algorithm

Input: A query node q to drive computation.
 A set of entities \mathcal{E} each with one mention m .
 A positive integer *samples*.
 A function \mathcal{I} that samples from nodes entities according to its influence on a mention.

Output: A set of resolved entities \mathcal{E}' .

```

1:  $\mathcal{E}' \leftarrow \mathcal{E} \cup q$ 
2: while samples -- > 0 do
3:    $m_1 \leftarrow \mathcal{I}(\mathcal{E}', q)$ 
4:    $m_2 \leftarrow \mathcal{I}(\mathcal{E}', q)$ 
5:    $\mathcal{E}'' \leftarrow \text{MOVE}(\mathcal{E}', m_1, m_2.\text{entity})$ 
6:   if  $\text{SCORE}(\mathcal{E}') < \text{SCORE}(\mathcal{E}'')$  then
7:      $\mathcal{E}' \leftarrow \mathcal{E}''$ 
8:   end if
9: end while
   return  $\mathcal{E}'$ 

```

For each iteration, the algorithm selects mentions using the influence function (line 3 and line 4). Then, one mention m_1 is moved into the entity of m_2 . Mentions m_1 and m_2 have a higher probability of being coreferent and therefore a higher probability of a merge occurring in the query entity compared to random selections as in

Algorithm 1. As a corollary, the influence sampling property creates many small entities that are similar to the query entity.

During query-proportional sampling more entities that are similar to the query node are created. Some of the mentions created in intermediate entities during query-proportional sampling will move to the query entity. This is a big advantage when performing entity-to-entity merges (as opposed to mention to entity merges). In this chapter, we do not investigate this extension to the algorithm.

5.3 Hybrid ER The best of both the target-fixed and query-proportional algorithms can be combined to create a hybrid algorithm. Like the target-fixed algorithm, the hybrid method aggressively fixes the target as the query entity. The hybrid method also chooses its source node using the influence function in the same manner as the query-proportional algorithm.

Algorithm 7 shows the hybrid algorithm using the repel method. With probability τ_α the algorithm chooses a mention using the repel method (\mathcal{I}) and moves it to an entity that is not the query node. This is the opposite of merging a node into the query entity. Pseudocode is listed on lines 3 to line 5.

Algorithm 7 Hybrid-Repel algorithm

Input: A set of entities \mathcal{E} , where one contains all the mentions m and the others are empty.
A positive integer *samples*.
A query node q .
A function \mathcal{I} that samples from nodes entities according to its influence on a mention.

Output: A set of resolved entities \mathcal{E}' .

```

1:  $\mathcal{E}' \leftarrow \mathcal{E} \cup q$ 
2: while samples -- > 0 do
3:   if RANDOM() <  $\tau_\alpha$  then
4:      $m \leftarrow \mathcal{I}(\mathcal{E}', q)$ 
5:      $e_i \leftarrow \{e \mid \exists e, e \in \mathcal{E}', e \neq q.\text{entity}\}$ 
6:   else
7:      $e_i \sim_u \mathcal{E}'$ 
8:      $e_j \leftarrow \{e \mid \exists e, e \in \mathcal{E}', e \neq e_i\}$ 
9:      $m \sim_u e_j$ 
10:  end if
11:   $\mathcal{E}'' \leftarrow \text{MOVE}(\mathcal{E}', m, e_i)$ 
12:  if SCORE( $\mathcal{E}'$ ) < SCORE( $\mathcal{E}''$ ) then
13:     $\mathcal{E}' \leftarrow \mathcal{E}''$ 
14:  end if
15: end while
    return  $\mathcal{E}'$ 

```

5.4 Implementation Details The previous algorithms described single process sampling over the set of mentions. The multi-query methods are modeled for several interwoven sequential single-node ER processes. In this section, we describe our implementation of the hybrid algorithm over a parallel database management system.

An independent Vose structure (\mathcal{Z} , § 5.1) is created for each query node in the query set. The creation of the Vose structure query nodes is parallelized. When the number of query nodes increases the Vose structures demand more memory from the system. Each Vose structure contains array of type double precision and unsigned int. The space for the structure is $O(|Q| \cdot |M|)$ where $|Q|$ is the number of query nodes in the query and $|M|$ is the number of mentions in the corpus. The Vose structure is accessed over every sample and needs to be in memory. To increase scalability, one could store the full sets of precomputed samples and serialize the Vose structures to disk but that is not explored here [19].

Sampling over the query nodes for each algorithm can also be perform in parallel. In our method, a thread

selects a query node using a random schedule as described in Section 4.3. The system will use the Vose structure associated with the query node to set up a proposal move. The system attempts to obtain a locks for both entities involved in the proposal. If the system is unable to obtain a lock on either of the two entities the system will back out and resample new entities. When the number of query nodes is small the query-driven algorithms experience lot of contention at the entities containing the query nodes. In these circumstances, the system will back out and either restart the proposal process or attempt a baseline proposal. This avoids waiting for locked entities and keeps the sampling process active. In Section 6.6 we demonstrate the parallel hybrid method over a large data set.

5.5 Algorithms Summary Discussion Algorithms 5, 6 and 7 are modifications of proposal jumps found in the baseline Algorithm 1. Table 3 describes the proposal process for each algorithm by its preferred jump method.

Table 3: Summary of algorithms and their most common methods for proposal jumps

	source	target
Baseline	random	random
Target-Fixed	random	fixed
Query-Proportional	proportional	proportional
Hybrid	proportional	fixed

The target-fixed algorithm builds the query entity by aggressively proposing random samples to merge into the query entity. The query-proportional algorithm uses an influence function to ensure its samples are mostly related to the query node. The hybrid algorithm mixes the aggressiveness of the target-fixed with the intelligent selecting of the source node found in the query proportional method.

After choosing the correct algorithm, a user needs to have a well trained model with several features. An advantage of using query-proportional techniques, because so little sampling is required, is that we can interactively test query accuracy. We can and also add context or keywords that were discovered from a previous run of the algorithm. This interactive querying workflow will help improve accuracy, which we experimentally verify in Section 6.

Parameter settings The algorithm takes several parameters that affect performance. While not studied in this chapter, parameter settings are robust to change making parameter selection simple. The first is the number of proposals (*samples*). This number can be a function on the size of the data set. Each query node should have the opportunity to be merged into an entity more than once.

The value τ_α is between $[0.0, 1.0]$ and represents how often to perform the main type of sampling. This value should be set to a high value, 0.9 for accept algorithms. With probability $1 - \tau_\alpha$ the algorithms back-off to random samples to improve mixing. This value is lowered to counter some of the aggression, particularly in Algorithm 5. The parallel experiments use a $\tau_\alpha = 1$ and back out when there is contention in the threads.

In statistics, a negative binomial function is used to model the number of trials it takes for an event to be a success. We can also use a negative binomial function as a decay function for the output of the influence function. We use this function because we want values that are most similar (lowest score) to be sampled more often. We set the r value, or number of failures for the negative binomial function to 1. We set the p value, or the probability of each success to a value close to 0.05.

In the multi-query ER algorithms we run inference for K steps before we look to change the query entity. In our experiments we choose a K of 500 and an increasing value from two to 100 thousand in the parallel experiments.

6 Query-Driven Entity Resolution Experiments

In this section, we describe the implementation details, the data sets and our experimental setup. Next, we discuss our hypotheses and four corresponding experiments. We then finish with a discussion of the results.

Implementation We developed the algorithms described in Section 4 in Scala 2.9.1 using the Factorie package. Factorie is a toolkit for building imperatively defined factor graphs [27]. This framework allows a templated definition of the factor graoh to avoid fully materializing the structure. The training algorithms are

also developed using Factorie. The algorithms for canopy building and approximate string matching are developed as inside of PostgreSQL 9.1 and Greenplum 4.1 using SQL, PL/pgSQL and PL/Python. Inference is performed in-memory on an Intel Core i7 processors with 3.2GHz, 8 cores and 12GB of RAM. The approximate string matching on Greenplum is performed on a AMD Opteron 6272 32-core machine with 64 GB.

The parallel experiments were developed entirely in a parallel database, DataPath [3]. DataPath is installed on a 48-core machine with 256 GBs.

Data sets. The experiments use three data sets, the first is the English newswire articles from the Gigaword Corpus, we refer to this as the NYT Corpus [15]. The second is a smaller but fully-labeled Rexa data set.³ Because it is fully-labeled it allows us to run the more detailed micro benchmarks. The NYT corpus contains 1,655,279 articles and 29,866,129 paragraphs from the years 1994 to 2006. We extracted a total of 71,433,375 mentions using the natural language toolkit named entity extraction parser [7]. Additionally, we compute general statistics about the corpus including the term and document frequency and tf-idf scores for all terms. We manually labeled mentions for each query over the NYT data set.

The second data set, Rexa, is citation data from a publication search engine named Rexa. This data set contains 2454 citations and 9399 authors of which 1972 are labeled. We perform experiments on the Rexa corpus because it is fully labeled unlike the NYT Corpus. The Rexa corpus is smaller in total size but it has average sized canopies.

The third data set is the Wikilinks Corpus [34] largest labeled corpus for entity resolution that we could find at the time of development. It contains 40 million mentions and 3 million entities that were extracted from the web and truthed based on web anchor links to Wikipedia pages. We loaded a million mentions onto DataPath to demonstrate the parallel capabilities.

6.1 Experiment Setup

Table 4 lists the features and the weights for each feature.

Features. Features that look for similarity between mention nodes are called affinity features and they are given positive weights. Features that look for dissimilarity between mentions nodes are called repulsion factors and they are given negative weights. We implement three classes of features: pairwise token features, pairwise context features and entity-wide features. Pairwise features directly compare tokens strings on attributes such as equality or matching substrings. Context features compare the information surrounding the mention. We can look at the surrounding sentence, paragraph, document or user specified keywords. The query nodes are extracted from text and contain a proper document context. With this context, we use a tf-idf weighted cosine similarity score to compare the context of each mention token. Finally, entity-wide features use all mentions inside an entity cluster to make a decision. An example entity-wide feature counts the matching mention strings between two entities.

Models. Features on the NYT and Wikilinks data sets were manually tuned and the features for the Rexa data set were trained using sample rank [41] with confidence weighted updates. We manually tune some of the weights in the NYT corpus to make up for the lack the complete training data. The models can be graphically represented as the models in Figure 4.

Evaluation metrics. Convergence of MCMC algorithms is difficult to measure as describe in a review by Cowles and Carlin [10]. We estimate the convergence progress by calculating the $f1$ score of the query node’s entity ($f1_q$). We create this new measure because we are primarily concerned with the query entity. Other measures include B^3 for entity resolution and several others for general MCMC models [4, 10].

The query-specific $f1$ score is the harmonic mean of the query-specific recall R_q and query-specific precision P_q . To accurately determine the P_q and R_q of each query in this experiment we label each correct query node. Query-specific precision is defined as $P_q = \frac{|\{\text{relevant}(\mathcal{M})\} \cap \{\text{retrieved}(\mathcal{M})\}|}{|\{\text{retrieved}(\mathcal{M})\}|}$ and query-specific recall $R_q = \frac{|\{\text{relevant}(\mathcal{M})\} \cap \{\text{retrieved}(\mathcal{M})\}|}{|\{\text{relevant}(\mathcal{M})\}|}$. The $f1$ score for the query node’s entity q is defined as:

$$f1_q = 2 \frac{R_q P_q}{R_q + P_q}.$$

The $f1_q$ score is a good indicator of entity and answer quality. For multi-query experiments we calculate the average $f1_q$ scores for each query node. The run of each non-parallel algorithm is averaged over 3 to 10 runs.

³<http://cs.neiu.edu/~culotta/data/rexa.html>

Table 4: Features used on the NYT Corpus. The first set of features are token specific features, the middle set are between pairs of mentions and the bottom set are entity wide features.

Feature name	Score ⁺	Score ⁻	Feature type
Equal mention Strings	+20	-15	Token Specific
Equal first character	+5		Token Specific
Equal second character	+3		Token Specific
Equal second character	0		
Unequal mention Strings		-15	Token Specific
Unequal first character	0		
Unequal second character	0		
Unequal second character	0		
Equal substrings	+30	-150	Token Specific
Unequal substrings		-150	Token Specific
Equal string lengths	+10		Token Specific
Matching first term	+90	-3	Token Specific
No matching first term		-3	Token Specific
Similarity ≥ 0.99	+120		Pairwise
Similarity ≥ 0.90	+105		Pairwise
Similarity ≥ 0.80	+80		Pairwise
Similarity ≥ 0.70	+55		Pairwise
Similarity ≥ 0.60	+35		Pairwise
Similarity ≥ 0.50	+15		Pairwise
Similarity ≥ 0.40		-5	Pairwise
Similarity ≥ 0.30		-50	Pairwise
Similarity ≥ 0.20		-80	Pairwise
Similarity < 0.20		-100	Pairwise
Matching terms	+20		Pairwise
Token in context	+1		Pairwise
No matching keyword	+700	-10	Pairwise
Matching Keyword	+700		Pairwise
Keyword in token	+70		Pairwise
Extra Token		-500	Pairwise
Matching token in context	+10		Pairwise
Similar neighbor	+100	-5	Entity-wide
No Similar neighbor in entity		-5	Entity-wide
Matching document	+350	-15	Entity-wide
No Matching documents in entity		-15	Entity-wide

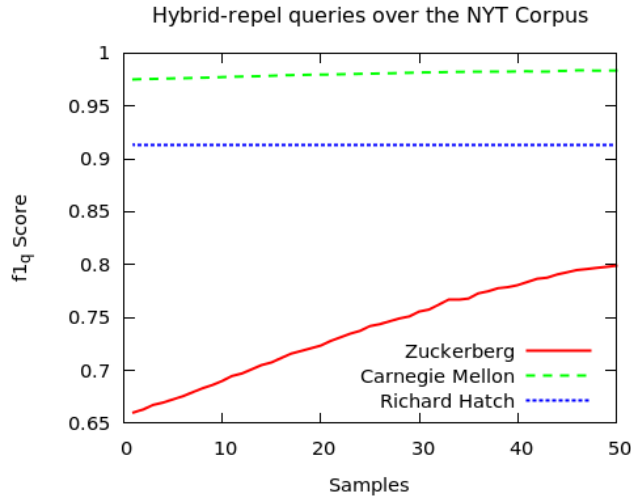


Figure 5: Hybrid-repel performance for the first 50 samples for three queries. Each result is averaged over 6 runs

6.2 Realtime Query-Driven ER Over NYT In this experiment we show that query-driven entity resolution techniques allow us to obtain near realtime⁴ results on large data sets such as the NYT corpus.

Figure 5 shows the $f1_q$ score of the hybrid ER algorithms with three single-query ER queries. The graph shows performance over the first 50 proposals. For example, the ‘Zuckerberg’ query could be expressed as shown in Algorithm 8.

Algorithm 8 Example ER query over the entity ‘zuckerberg’

```
SELECT * FROM
Mention m
WHERE coref_map(m.*, entity^p), ‘zuckerberg’, context).
```

Recall, a canopy is first generated using an approximate match over the mention set. We use the repel inference function and all the mentions are initialized in one large entity. The ‘Richard Hatch’ and ‘Carnegie Mellon’ queries start at an $f1_q$ score of .92 and .97, respectively. The ‘Zuckerberg’ query starts above .65 and improves to an $f1_q$ score over .8.

These experiments show the repel method removing mismatches from the query entity. The co-influence function is used to quickly identify the mentions that do not belong in the entity and they are proposed to be removed. When a hybrid move is proposed, a mention from the large entity moved from a large entity group to a new, possibly empty, entity. This method relies on the good repulsion features and correct weights.

In Table 5 we show the performance of three queries. In addition to the query token we add four columns: blocking time in seconds, canopy size, inference time in seconds and the total compute time. Total time is the complete time taken by each run, this includes building of the influence data structure and result writing. The values in Table 5 show that fast performance of query-driven ER over a large database of mentions.

6.3 Single-query ER In this experiment we show a performance comparison between the single-query algorithms summarized in Sections 4 and 5. We run the query-driven algorithms over queries with different selectivity levels and show the accuracy over time. Each algorithm uses the attract method, so each mention in the canopy starts in its own entity.

⁴We define realtime as only contributing a small or no time loss when this process a part of an external execution pipeline such as an information extraction pipeline.

Table 5: The performance of the hybrid-repel ER algorithm for queries over the NYT corpus for the first 50 samples. Total time includes the time to build the $\bar{\mathcal{I}}$ data structure and result output. The NYT Corpus contains over 71 million mentions, a large amount for the entity resolution problems.

Query	Blocking	Mentions	Inference	Total time
Zuckerberg	24.4 s	103	2 s	37 s
Richard Hatch	28.3 s	226	18.5 s	59 s
Carnegie Mellon	25.9 s	1302	68 s	124 s

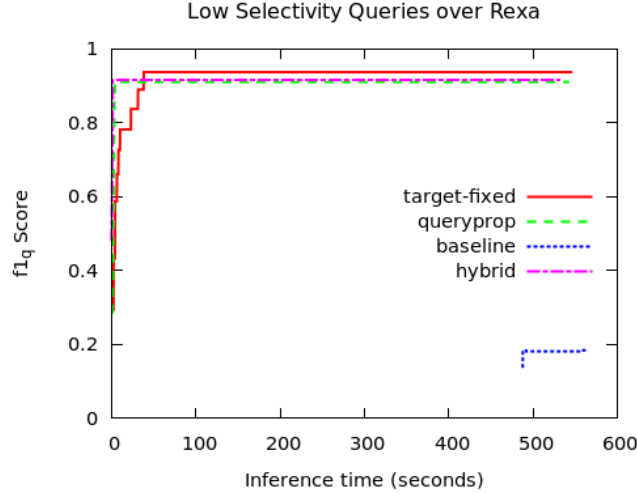


Figure 6: A comparison of single-query algorithms on a query with selectivity of 11

Figure 6 shows the run time of all four algorithms on the Rexa data set with the query ‘Nemo Semret’, an author with a selectivity of 11. The performance for the baseline entity resolution does not get a correct proposal until about 500 seconds. The baseline algorithm takes a long time to accept the first proposal because it is randomly trying to insert mentions into an existing entity. Target-fixed immediately begins to make correct proposals. Hybrid and query-proportional have the best performance and resolve the entity almost instantaneously. The hybrid chooses the most likely nodes to merge into the query entity. As the first couple of proposals are correct merges, hybrid quickly converges to a high accuracy. Due to imperfect features, among the 10 averaged runs a few runs get stuck at local optimum and causing suboptimal results.

Figure 7 shows the run time of four algorithms for query node id ‘A. A. Lazar’ with selectivity of 46. The baseline algorithm progresses the slowest. The hybrid algorithm quickly reaches a perfect $f1_q$ score. Query-proportional algorithm lags slightly behind the hybrid method but still reaches a perfect value. The target-fixed algorithm gradually increases to a perfect $f1_q$ score about 60 seconds after hybrid and query-proportional.

Figure 8 shows the run time of four algorithms with a query ‘Michael Jordan’ of selectivity 130. The baseline slowly increases over the 100 seconds. The hybrid algorithm again quickly achieves a perfect $f1_q$ score followed by query-proportional and then target-fixed. The time gap between each of the algorithms increases with the increase in selectivity, hybrid achieves the best performance.

We look deeper at how selectivity affects the rate of convergence. In Figure 9 we show the time it takes for each algorithm to reach an $f1_q$ score of 0.95 over increasing selectivity. We choose five query nodes of increasing selectivity but with the same canopy sizes. The hybrid algorithm runtime increased with the increase in selectivity but only slightly steeper than constant. Target-fixed increased for the first three queries but did not last more than 50 seconds. Query-proportional has only a slight increase in time till convergence for the first three queries. The highest two selectivity queries are expensive for query-proportional and we observe an exponential increase in runtime. These results are consistent with the exponentially large increase in the number of random comparisons needed to find a match for a query entity. The query-proportional algorithm does not focus on the query entity

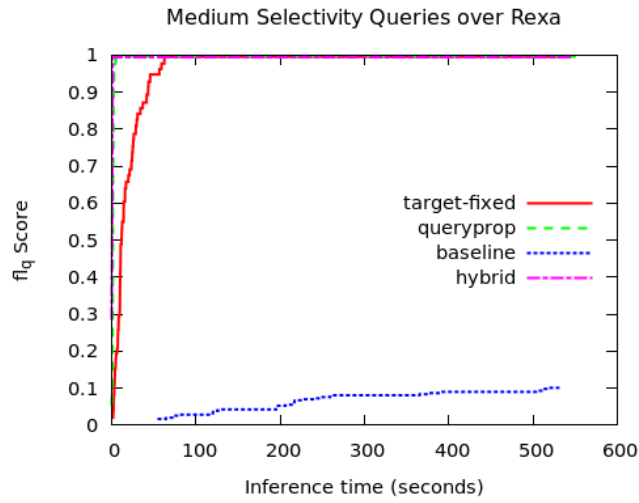


Figure 7: A comparison of single-query algorithms with a query node of selectivity 46

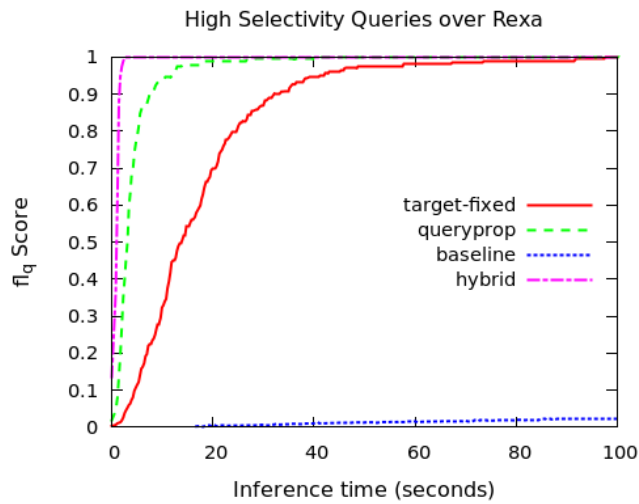


Figure 8: A comparison of selection-driven algorithms with a query node of selectivity 130

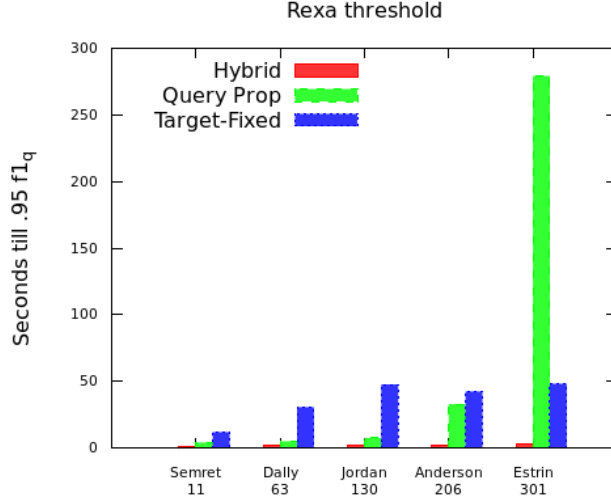


Figure 9: The time until an $f1_q$ score of 0.95 for five queries of increasing selectivities; averaged over three runs

as aggressively as target-fixed and hybrid algorithms. Recall that the target-fixed and the hybrid algorithm focus on moving correct nodes into the query entity. Query-proportional selects candidate nodes using the influence function but it does not fix the target entity. With the target entity not fixed, the chance of correct node for the query entity decrease exponentially. This shows that selectivity of nodes affects the runtime performance of each algorithm. When performing join-driven ER it is important to take the relative selectivity of nodes into account for choosing best scheduling algorithms.

6.4 Multi-query ER In this experiment we study performance of our different scheduling algorithms for join-driven ER queries. We choose ten query nodes of different selectivity and run the join queries scheduling algorithms described in Section 4.3. Consider a table like the People table in Section 3 with selectivity $\{130, 63, 68, 7, 12, 12, 301, 11, 46\}$. The four algorithms, random, closest-first, farthest-first and selectivity-based are shown in Figure 10. The selectivity-based method out performs the other three algorithms in terms of convergence rate. The jumps in accuracy on the graph correspond to the scheduling algorithms choosing new query nodes and accepting new proposals. It has a high jump when it starts sampling the seventh, and highest selectivity nodes. The farthest-first algorithm rises the slowest out of the scheduling algorithms because it tries to stop sampling the high performing query entity and makes proposals for the slowest growing. Selectivity-based method performs well early because the high selectivity queries are sampled first. The high selectivity query makes up a large proportion of the total $f1_q$ score. The large jump in the random method is when it reaches the node with selectivity 301. Notice, closest-first reaches its peak $f1_q$ score the fastest because it tries to get the most out of every query node.

6.5 Context Levels In this experiment we aim to discover how different levels of context specified at query time can improve convergence time and overall accuracy. We take the zuckerberg query and the hybrid-repel algorithm and ran ER three times over three levels of context. Each mention in the graph contains a ‘paragraph’ level of context and we only alter the context of the query node. The ‘none’ context only activates token specific features, any context features involving the query node are zeroed out. The ‘paragraph’ level context is the default context from the NYT corpus and the ‘document’ level context extends context to the entire news article. Additionally, we add specific keywords from Mark Zuckerberg’s DBpedia page to the ‘document’ and ‘paragraph’ context levels. We show the performance using the repel method in Figure 11.

Adding specific keywords that activate the keyword features are the most effective methods for increasing the accuracy of query-driven ER. Query-driven methods allow a user to observe the results and add or remove keywords specific queries to improve the accuracy. This type of iterative improvement workflow is not feasible with batch methods.

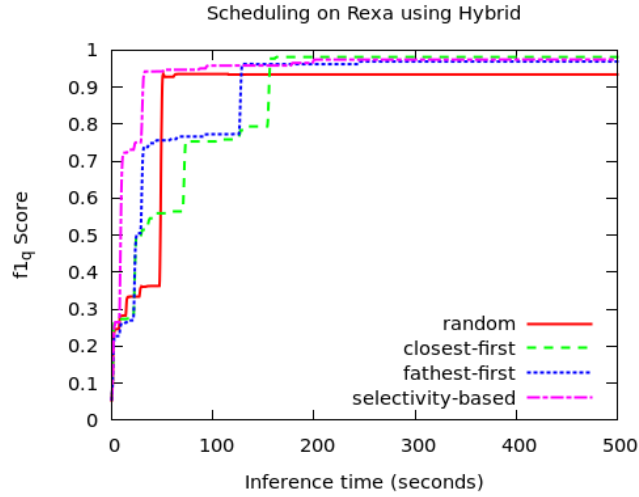


Figure 10: The progress of the hybrid algorithm across for multiple query nodes using difference scheduling algorithms. Each result is averaged over three runs

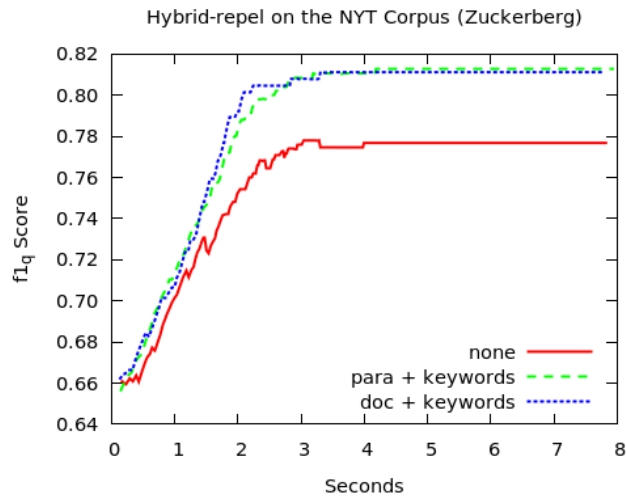


Figure 11: The performance of zuckerberg query with difference levels of context. Each result is averaged over 6 runs

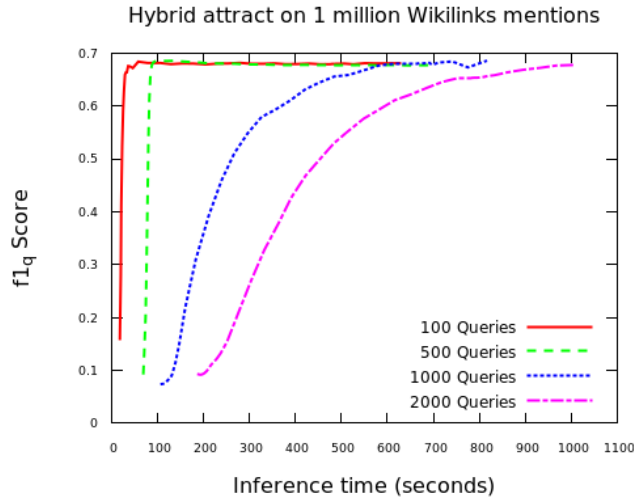


Figure 12: Hybrid-attract algorithm with random queries run over the Wikilinks corpus. Each plot starts after the Vose structures are constructed

6.6 Parallel Hybrid ER In this experiment has two objectives, first how does the hybrid algorithm perform in a canopy size of 1 million queries and what is the effect of increasing the number of queries nodes. In Figure 12 the Hybrid algorithm is able to resolve entities in a short amount of time. The creation time of the Vose structure is about linear in the number of queries. The trend in the graph is that as the ratio of queries to entities increases the performance benefit of the hybrid-attract method decreases. With more query nodes the construction time increases and the benefits of the algorithm decrease and become no better than the baseline method.

Experiment Summary Each of the query-driven methods outperform the baseline methods in terms of runtime while not losing out on accuracy. Across different data set sizes hybrid algorithms have the most consistent performance. If a system has a quality blocking function then it is better to use the co-influence entity resolution method. With multiple query nodes, selectivity-based is the most consistent performing algorithm. More accurate estimation of MCMC convergence performance could allow the dynamic scheduling algorithms closest-first and farthest-first to achieve higher accuracy. The more contextual information that can be added to query nodes at query time causes higher accuracy of the entity resolution algorithms. Parallel query-driven sampling is an effective way to get speed up in an ER data set when the ratio of mentions to entities is low.

7 Query-Driven Entity Resolution Related Work

This chapter is related to work in several areas. In this section we describe a selection of the literature that we found most relevant to different parts of the Query-Driven ER task.

Entity Resolution The state-of-the-art method for entity resolution employs collective classification. Instead of purely pairwise decisions, collective classification methods consider group relationships when making clustering determinations. In a recent tutorial [14], collected classification methods were grouped into three categories: non-probabilistic [5, 12, 20], probabilistic [8, 13, 23, 28, 29, 35] and hybrid approaches [2, 31]. A relevant challenge proposed for entity resolution research by the tutorial is how to efficiently perform entity resolution when a query is involved. This chapter seeks to address this issue.

Entity resolution is generally an expensive, offline batch process. Bhattacharya and Getoor proposed a method for query-time entity resolution [6]. This method performs inference by starting with a query node and performing ‘expand and resolve’ to resolve entities through resolution of attributes and expansion of hyper-edges. Unfortunately, hyper-edges between records are not always explicit in data sets. This chapter does not assume the presence of any link in the corpus, each entity or mentions are independently defined, which is the case for most applications.

A recent paper by Altwaijry, Kalashnikov and Mehrotra [1] has a similar motivation of using SQL queries to drive entity resolution. That work focuses using predicates in the query to drive computation while this work uses

example queries to drive computation. Both techniques are complementary and combining the two by updating the edge-picking policy described in their paper using our approach makes for interesting method of optimizing the entity resolution process.

The term query-driven appears in this chapter and has appeared in others across literature with different meanings [16]. Our definition of a query node is an example item, mention, from a data set. A query in Altwaijry et al. [1] are the predicates in an SQL statement. Query-driven in Grant et al. [16] is the SQL queries used to drive analytics.

It is becoming increasingly normal to work with data sets of extremely large size, in response researchers have studied streaming and distributed processing. Rao, McNamee and Dredze describes an approach for streaming entity resolution [30]. This approach is fast and approximates entries in an LRU queue of clustered entity chains. We apply these techniques to static data set and do not yet handle streams of data. Singh, Subramanya, Pereira, and McCallum propose a technique for ER where entities are resolved in parallel blocks and then redistributed and resolved again in new blocks [33]. This parallel distribution method makes large-scale entity resolution tractable. In this chapter, we perform analysis on a similar scale data set but we show that great performance gains can be achieved when a query is specified.

Query specific sampling Recently, several researchers have explored the idea of focusing sampling of graphical models to speed up inference. Below we discuss the three approaches that use sampling to speed up ER over graphical models.

Query-Aware MCMC [40] found that when performing a query over a graphical model the cost of not sampling a node is exactly the nodes influence on the query node. This enables us to ignore some nodes that have low influence over the query node and incur a small amount of error. This influence score can be calculated as the mutual information between two nodes. The authors compare estimation techniques of the intractable mutual information score, this is called the influence trail score. Because ER has a fixed pairwise model, we can use the theory from this work and specialized data structures to gain performance when query-driven sampling.

Type-based MCMC is a method of sampling groups of nodes of with the same attribute to increase the progress towards convergence [24]. This approach works well when feature sets can be tractably counted and grouped. If query nodes are introduced it is not clear how one may focuses type-based sampling.

Other researchers have explored using belief propagation with queries to approximate marginal of factor graphs [9]. However, the entity resolution graph is cyclic and highly connected. MCMC scales with large real world models better than loopy belief propagation [40].

8 Query-Driven Entity Resolution Summary

In this chapter, I propose new approaches for accelerating large-scale entity resolution in the common case that the user is interested in one or a watch list of entities. These techniques can be integrated into existing data processing pipelines or used as a tool for exploratory data analysis. We showed three single-node ER algorithms and three scheduling algorithms for multi-query ER and show experimentally how their runtime performance is several orders of magnitude better than the baseline.

References

- [1] H. Altwaijry, D. V. Kalashnikov, and S. Mehrotra. Query-driven approach to entity resolution. *Proceedings of the VLDB Endowment*, 6(14), 2013.
- [2] A. Arasu, R. Christopher, and D. Suciu. Large-scale deduplication with constraints using dedupalog. In *International Conference on Data Engineering*, pages 952–963. IEEE, 2009.
- [3] S. Arumugam, A. Dobra, C. M. Jermaine, N. Pansare, and L. Perez. The datapath system: A data-centric analytic processing engine for large data warehouses. In *Proc. of the 2010 ACM SIGMOD*, pages 519–530, NY, USA, 2010. ACM.
- [4] A. Bagga and B. Baldwin. Entity-based cross-document coreferencing using the vector space model. In *17th ACL*, pages 79–85. ACL, 1998.
- [5] I. Bhattacharya and L. Getoor. Collective entity resolution in relational data. *ACM Trans. KDD*, 1(1), Mar. 2007.
- [6] I. Bhattacharya, L. Getoor, and L. Licamele. Query-time entity resolution. In *Proc. 12th ACM SIGKDD*, KDD ’06, pages 529–534, NY, USA, 2006.
- [7] S. Bird, E. Loper, and E. Klein. Natural language processing with python. In *Natural Language Processing with Python*. O’Reilly Media Inc, 2009.

- [8] M. Bröcheler, L. Mihalkova, and L. Getoor. Probabilistic similarity logic. In P. Grünwald and P. Spirtes, editors, *UAI*, pages 73–82. AUAI Press, 2010.
- [9] A. Checheta and C. Guestrin. Focused belief propagation for query-specific inference. In *AISTATS*, May 2010.
- [10] M. Cowles and B. Carlin. Markov chain monte carlo convergence diagnostics: a comparative review. *Journal of AmStat*, 91(434):883–904, 1996.
- [11] A. Das Sarma, A. Jain, A. Machanavajjhala, and P. Bohannon. An automatic blocking mechanism for large-scale de-duplication tasks. In *Proceedings of the 21st ACM CIKM*, pages 1055–1064. ACM, 2012.
- [12] X. Dong, A. Halevy, and J. Madhavan. Reference reconciliation in complex information spaces. In *Proceedings of the 2005 ACM SIGMOD international conference on Management of data*, SIGMOD ’05, pages 85–96, New York, NY, USA, 2005. ACM.
- [13] I. Getoor. A latent dirichlet model for unsupervised entity resolution. In *Proceedings of the 6th SIAM International Conference on Data Mining*, volume 124, page 47. Society for Industrial Mathematics, 2006.
- [14] L. Getoor and A. Machanavajjhala. Entity resolution: Theory, practice & open challenges. In *Proceedings of the 38rd VLDB*, VLDB ’12. VLDB Endowment, 2012.
- [15] D. Graff. Ldc2007t07: English gigaword corpus, 2007.
- [16] C. E. Grant, J.-d. Gumbs, K. Li, D. Z. Wang, and G. Chitouras. Madden: query-driven statistical text analytics. In *Proceedings of the 21st ACM international conference on Information and knowledge management*, CIKM ’12, pages 2740–2742, New York, NY, USA, 2012. ACM.
- [17] L. Gravano, P. Ipeirotis, H. Jagadish, N. Koudas, S. Muthukrishnan, L. Pietarinen, and D. Srivastava. Using q-grams in a dbms for approximate string processing. *IEEE Data Engineering Bulletin*, 24(4):28–34, 2001.
- [18] J. M. Hellerstein, C. Ré, F. Schoppmann, D. Z. Wang, E. Fratkin, A. Gorajek, K. S. Ng, C. Welton, X. Feng, K. Li, and A. Kumar. The madlib analytics library: or mad skills, the sql. *Proceedings of the VLDB Endowment*, 5(12):1700–1711, Aug. 2012.
- [19] R. Jampani, F. Xu, M. Wu, L. L. Perez, C. Jermaine, and P. J. Haas. Mcdm: A monte carlo approach to managing uncertain data. In *Proc. of the 2008 ACM SIGMOD*, pages 687–700, NY, USA, 2008. ACM.
- [20] D. V. Kalashnikov and S. Mehrotra. Domain-independent data cleaning via analysis of entity-relationship graph. *ACM Trans. Database Syst.*, 31(2):716–767, June 2006.
- [21] D. Koller and N. Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- [22] K. Li, C. Grant, D. Z. Wang, S. Khatir, and G. Chitouras. Gptext: Greenplum parallel statistical text analysis framework. In *Proceedings of the Second Workshop on Data Analytics in the Cloud*, pages 31–35. ACM, 2013.
- [23] X. Li, P. Morie, and D. Roth. Identification and tracing of ambiguous names: Discriminative and generative approaches. In *Proceedings of the National Conference on Artificial Intelligence*, pages 419–424. Menlo Park, CA; Cambridge, MA; London; AAAI Press; MIT Press; 1999, 2004.
- [24] P. Liang, M. I. Jordan, and D. Klein. Type-based mcmc. In *Human Language Technologies: The 2010 NAACL, HLT ’10*, pages 573–581, Stroudsburg, PA, USA, 2010. ACL.
- [25] J. Madhavan, S. Jeffery, S. Cohen, X. Dong, D. Ko, C. Yu, and A. Halevy. Web-scale data integration: You can only afford to pay as you go. In *Proceedings of CIDR*, pages 342–350, 2007.
- [26] A. McCallum, K. Nigam, and L. Ungar. Efficient clustering of high-dimensional data sets with application to reference matching. In *Proc. of the 6th SIGKDD*, pages 169–178, 2000.
- [27] A. McCallum, K. Schultz, and S. Singh. FACTORIE: Probabilistic programming via imperatively defined factor graphs. In *NIPS*, pages 1426–1427, 2009.
- [28] A. McCallum and B. Wellner. Conditional Models of Identity Uncertainty with Application to Noun Coreference. In *NIPS*, 2004.
- [29] H. Pasula, B. Marthi, B. Milch, S. J. Russell, and I. Shpitser. Identity Uncertainty and Citation Matching. In *Neural Information Processing Systems*, pages 1401–1408, 2002.
- [30] D. Rao, P. McNamee, and M. Dredze. Streaming cross document entity coreference resolution. In *Proceedings of the 23rd International Conference on Computational Linguistics: Posters*, pages 1050–1058. Association for Computational Linguistics, 2010.
- [31] W. Shen, X. Li, and A. Doan. Constraint-based entity matching. In *Proceedings of the National Conference on Artificial Intelligence*, volume 20, page 862. Menlo Park, CA; Cambridge, MA; London; AAAI Press; MIT Press; 1999, 2005.
- [32] L. Shu, A. Chen, M. Xiong, and W. Meng. Efficient spectral neighborhood blocking for entity resolution. In *2011 IEEE 27th ICDE*, pages 1067–1078, april 2011.
- [33] S. Singh, A. Subramanya, F. Pereira, and A. McCallum. Large-scale cross-document coreference using distributed inference and hierarchical models. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*, pages 793–803. Association for Computational Linguistics, 2011.
- [34] S. Singh, A. Subramanya, F. Pereira, and A. McCallum. Wikilinks: A large-scale cross-document coreference corpus

labeled via links to Wikipedia. Technical Report UM-CS-2012-015, 2012.

- [35] P. Singla and P. Domingos. Entity Resolution with Markov Logic. In *IEEE International Conference on Data Mining*, pages 572–582. IEEE, 2006.
- [36] M. D. Vose. A linear algorithm for generating random numbers with a given distribution. *IEEE Trans. Softw. Eng.*, 17(9):972–975, Sept. 1991.
- [37] S. E. Whang, D. Menestrina, G. Koutrika, M. Theobald, and H. Garcia-Molina. Entity resolution with iterative blocking. In *2009 ACM SIGMOD*, pages 219–232. ACM, 2009.
- [38] M. Wick, A. McCallum, and G. Miklau. Scalable probabilistic databases with factor graphs and mcmc. *Proc. VLDB Endow.*, 3(1-2):794–804, Sept. 2010.
- [39] M. Wick, S. Singh, and A. McCallum. A discriminative hierarchical model for fast coreference at large scale. In *Proceedings of the 50th ACL*, ACL ’12, pages 379–388, 2012.
- [40] M. L. Wick and A. McCallum. Query-aware mcmc. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K. Weinberger, editors, *Advances in NIPS 24*, pages 2564–2572, 2011.
- [41] M. L. Wick, K. Rohanimanesh, K. Bellare, A. Culotta, A. McCallum, and A. McCallum. Sample rank: Training factor graphs with atomic gradients. In *ICML*, pages 777–784, 2011.