

Web Scale Entity Resolution using Relational Evidence

Taesung Lee ^{†,‡,*} Zhongyuan Wang [‡] Haixun Wang [‡] Seung-won Hwang [†]
[†]POSTECH
Korea, Republic of
{elca4u,swhwang}@postech.edu
[‡]Microsoft Research Asia
Beijing, China
{zhy.wang,haixunw}@microsoft.com

ABSTRACT

Entity resolution has been extensively studied. Many approaches have been proposed, including using machine learning techniques to derive domain-specific lexical similarity measures, or rank entities' attributes by their discriminative power, etc. In this paper, we study the problem in the setting of matching two web scale taxonomies. Besides the scale, we address the challenge that the taxonomies may not contain enough context (such as attributes) for entity resolution, and traditional lexical similarity measures result in many false positive matches. To tackle this new task, we explore negative evidence in the structure of the taxonomy, as well as in external data sources such as the web. To integrate positive and negative evidence, we formulate the entity resolution problem as a problem of finding optimal multi-way cuts in a graph. We analyze the complexity of the problem, and propose a Monte Carlo algorithm for finding greedy cuts. We conduct extensive experiments that demonstrate the advantage of our approach.

1. INTRODUCTION

One of the tasks in data cleaning and data integration is entity resolution. Recently, entity resolution arises as one of the biggest challenges in the area of ontology/taxonomy managing, and in particular, in ontology/taxonomy mapping.

As understanding data semantics becomes more and more important in many applications, researchers and industrial practitioners create all kinds of ontologies and taxonomies to manage semantics. Each taxonomy has its own scope and characteristics, and many also have significant overlap. It is important to create mappings between multiple taxonomies for two reasons. First, mappings improve understanding, even if the two taxonomies are in the same domain. Second, taxonomy construction is a costly process, and it saves time and money if an existing taxonomy can be used to enrich a new taxonomy, and vice versa.

We focus on two web-scale taxonomies, namely Freebase [8] and Probase [35], to study this problem. From their major features listed in Table 1, we can see they are unique in their own ways.

*The work is conducted while the author works at Microsoft Research Asia

Probase is automatically harvested from web data. Although both Probase and Freebase are big, Probase is unique in the sense that it has an extremely large number of categories (2 million in Probase vs. 12.7 thousand in Freebase). For example, in Probase, *Emerging Markets* is a category on its own, and it includes subcategories such as *BRIC*, and instances or entities such as *China*, *Mexico*, *Russia*, *India*, *Brazil*, etc. The goal of Probase is to cover as many concepts of worldly facts in the collective mind of human beings as possible, in the hope that it will enable better understanding of human communications (natural language).

	Freebase	Probase
how is it built?	manual	automatic
data model	deterministic	probabilistic
taxonomy topology	mostly tree	DAG
# of categories	12.7 thousand	2 million
# of entities	13 million	16 million
information about entity	rich	sparse
adoption	widely used	new

Table 1: Two unique taxonomies.

On the other hand, Freebase has richer information about many entities. For example, for someone like Barack Obama, Freebase has the date of his birth, names of his spouse and kids, information about his religion, political party, etc. Although Probase contains a list of attributes for each category (i.e., Probase knows birthday is an attribute of a person), there are not many attribute values. Thus, a big motivation for matching the two taxonomies is to enrich the content of Probase at the entity level, and to enrich the content of Freebase at the category level.

However, matching two web-scale taxonomies such as Freebase and Probase pose at least two big challenges:

- Scalability. Probase has 16 million entities and Freebase has 13 million entities. It is infeasible to perform pairwise comparison for entity resolution.
- Insufficient context for entity resolution. Since Probase lacks attribute values for many entities, we cannot decide whether “J. Doe” is the same person as “John Doe” by comparing their birthdays.

Overcoming the above two challenges is essential to effectively creating a mapping between two web-scale taxonomies. To address the first problem, instead of comparing every pair of entities, we first look at the categories they belong to in each taxonomy. The intuition is that we probably do not need to compare entities in the category of “rare plants” with entities in the category of “an-

imals.” However, we might need to associate the category of “endangered species” with the category of “animals.” The challenge is that among the $12.7 \text{ thousand} \times 2 \text{ million}$ pairs of categories between Freebase and Probase, how do we figure out which pairs need our attention?

To address the second problem of not having enough contextual information for entity resolution, we collect *relational evidence*, which is evidence derived from relationships among entities, rather than from entities’ content. Specifically, we collect two types of relational evidence. The first type is positive evidence. Lexical similarity is a source of non-relational positive evidence. However, lexical similarity is not sufficient, as for example, “New York City” and “the Big Apple” have no lexical similarity, but they refer to the same entity. In this paper, we collect positive evidence by focusing on entities’ relationships manifested in data sources such as Wikipedia or the world wide web.

Yet a more important type of relational evidence is negative evidence. If we have a sentence that says “... presidents such as *George W. Bush*, Bill Clinton, and *George H. W. Bush*”, then we may conclude that *George W. Bush* and *George H. W. Bush* refer to two different persons, because a well formed list probably contains no duplicates. We show that such kind of negative evidence can be found in many external sources, for example, in tables or lists on the web, or in a very well organized taxonomy (such as Freebase).

Finally, given both positive and negative evidence, we are facing the task of consolidating them to achieve entity resolution of high accuracy. For instance, knowing that x is similar to y , and y is similar to z (positive evidence), and also knowing that x is not z (negative evidence), how do we perform entity resolution for x , y , and z ? In this paper, we formulate the problem of entity resolution with positive/negative evidence as a multi-way graph cut problem, and we propose a greedy approach to effectively solve the problem.

The rest of the paper is organized as follows. Section 2 describes the taxonomies we are working with. In Section 3, we introduce the types of evidence we use for entity resolution. Section 4 presents our method of entity resolution and taxonomy matching. In Section 5, we report experimental results, and we conclude in Section 6.

2. TAXONOMIES

A taxonomy or an ontology provides a shared conceptualization of a domain. Recently, there is a lot of interest in using structured data to empower search or other applications. A general purpose taxonomy about worldly facts is indispensable in understanding the user intent, and much efforts are being devoted to composing and managing such taxonomies.

Freebase [8] is a taxonomy composed mainly by its community members. It is an online collection of structured data harvested from many sources, including individual ‘wiki’ contributions. Freebase aims to create a global resource which allows people (and machines) to access common information more effectively.

Compared with manually constructed taxonomies, taxonomies automatically generated from data have advantages in scale and costs. Probase [35] is a research prototype that aims at building a unified taxonomy of worldly facts from web data and search log data. Compared with Freebase, the Probase taxonomy is extremely rich. The core taxonomy alone (which is learned from 1.68 billion web pages and 2 years’ worth of Microsoft Bing’s search log) contains more than 2 million categories, while Freebase contains about 12.7 thousand categories. As categories in Probase correspond to concepts in our mental world, Probase is valuable to a wide range of applications, such as search [32], where there is a need to interpret users’ intent.

Probase contains many isA relationships that are harvested using

the so called Hearst linguistic patterns [24], that is, SUCH AS like patterns. For example, a sentence that contains “... artists such as Pablo Picasso ...” can be considered as an evidence for the claim that *Pablo Picasso* is an instance in the *artist* category. For each category, Probase also collects a large set of attributes that can be used to describe instances in the category. For instance, the *artist* category may contain such attributes as *name*, *age*, *nationality*, *genre*, *specialization*, etc. Furthermore, Probase contains many relationships among instances of different categories. Fig. 3 (in Appendix) shows an interface with which users can browse the Probase taxonomy, and we can see a category (politicians) has many super categories, sub categories, instances, and similar categories.

It is not difficult to see that there is a strong need to integrate a taxonomy like Freebase with a taxonomy like Probase. With the integration, Freebase will have more information about categories, allowing Freebase to understand human concepts better, and Probase will have more information for each instance, giving Probase more knowledge in inference. In this paper, we study the challenges in creating a mapping between such taxonomies under this setting.

3. RELATIONAL EVIDENCE

In this section, we discuss how to obtain and quantify relational evidence, and in Section 4 we discuss how to use the evidence for entity resolution. Previous work assumes that each entity comes in a context, for example, a text window where an entity appears, or a set of attributes (such as a person’s gender or birthday) that describe the entity. Such context information is then used as positive or negative evidence. However, the context may be insufficient or noisy (e.g., the text window around an entity may contain irrelevant information that confuses entity resolution). In our work, we focus on a large number of entities that come with little context. For instance, assume a list contains nothing but three names *George W. Bush*, *George H. W. Bush*, and *Dubya*, how to find out how many distinct entities it contains? In this section, we explore evidence beyond entities’ immediate context to solve this problem.

3.1 Negative Evidence

Let (x_i, y_i) denote the claim that x_i and y_i represent the same entity. Any evidence that supports the claim (x_i, y_i) is called positive evidence, and any evidence that rejects the claim (x_i, y_i) is called negative evidence.

Negative evidence is particularly important in entity resolution. For instance, string similarity may provide strong evidence that *George W. Bush* is likely *President Bush*, and *President Bush* is likely *George H. W. Bush*. Intuitively, we may conclude that *George W. Bush* is *George H. W. Bush*, unless there is negative evidence to break the transitivity.

The challenge is then, how to find negative evidence? Previous work is based on data content. For example, if two persons with similar names have different birthdays, then we can conclude that they cannot be the same individual, unless the data is wrong. However, in many cases (e.g., taxonomies such as Probase), we do not have sufficient content information for all entities.

We argue that although individual data items may not contain much content, the community formed by the related items may contain valuable clue to entity resolution. In our work, we derive negative evidence based on how the data is inter-connected internally (e.g., in the taxonomies we are studying) as well as externally (e.g., on the web).

The ‘Birds of a Feather’ Principle (BoF)

Here is our intuition: *Michael Jordan* the professor and *Michael Jordan* the basketball player may not have too many friends in com-

mon. In other words, unless two persons with similar names have many common friends (i.e., their friends also have similar names), it is not likely that the two names refer to the same individual.

In a taxonomy (such as Probase), a data element may not contain much content, but the connections among the elements can be extremely rich. We use the connections, or the structure formed by the data, as a most important source of evidence. Specifically, we consider two types of connections: i) the connection between elements and the category they belong to (e.g., *{Michael Jordan, Shaquille O'Neal}* and *basketball players*); and ii) the connection between attributes and the category they describe (e.g., *{genre, artist, producer}* and *album*). In a modern taxonomy such as Probase, categories, elements, and attributes are all objects, and they are interconnected into a graph. Thus, a category can be regarded as a graph community formed by elements and attributes. We derive important negative evidence from such communities:

Negative evidence from the community: Let t be a threshold, and c_1, c_2 be two categories. We consider $\text{sim}(c_1, c_2) \leq t$ as negative evidence for any claim (x, y) where $x \in c_1$, and $y \in c_2$.

Here, $\text{sim}(c_1, c_2)$ measures the similarity between two categories. We define $\text{sim}(c_1, c_2)$ by linear combination as follows:

$$\text{sim}(c_1, c_2) = \lambda \cdot f(E_{c_1}, E_{c_2}) + (1 - \lambda) \cdot f(A_{c_1}, A_{c_2}) \quad (1)$$

where E_c and A_c denote the set of elements and attributes in category c respectively, λ is a parameter that balances the importance between elements and attributes, and f is a set similarity function based on Jaccard distance.

The ‘Clean data has no duplicates’ Principle (CnD)

If a list is well formed, then it probably does not contain any duplicates. Freebase is manually created and maintained. Thus, we can be relatively certain that each element in a Freebase category is a unique element in that category. This gives us negative evidence for entity resolution. For example, given that *George W. Bush* and *George H. W. Bush* both appear in the category of *US Presidents*, we can conclude that the two very similar names cannot be referring to the same individual.

Probase, on the other hand, is automatically created and maintained. Thus, a category may contain duplicates. For example, *Bill Clinton* and *William J. Clinton* may both appear in the category of *US Presidents*. Still, we can derive negative evidence from Probase. As we mentioned, Probase derives the isA relationship from Hearst patterns. For example, from the sentence “US Presidents such as George W. Bush, Bill Clinton, George H. W. Bush,” Probase concludes that *George W. Bush*, *Bill Clinton*, *George H. W. Bush* are US Presidents. But given that the three names appear in the same sentence, we can conclude that the three names represent three individuals. In other words, there is negative evidence for the claim (*George H. W. Bush*, *George W. Bush*), but there is likely no negative evidence for (*Bill Clinton*, *William J. Clinton*).

Besides Freebase and Probase, we also derive negative evidence from the web using the same argument. One source of evidence comes from Wikipedia. Wikipedia contains many lists such as a list of mountains. Furthermore, the lists are well formed and easily identifiable: they all have a title in the form of “list of ***.” We extract entity names from the list, and assume mentions in the list will represent different entities. Furthermore, Wikipedia has structured tables. Usually, a table has an entity column and multiple attribute columns. In the entity column, we can get entity names and assume there are no duplicates. In theory, we can apply the same reasoning to data on the web. However, web data is often very noisy, which

reduces the quality of the negative evidence.

3.2 Positive Evidence

Each piece of positive evidence is associated with a weight in the range of $[0, 1]$, which indicates how strong the evidence is. For instance, string similarity can be used as a source of evidence. Table 2 gives some examples, where we use Jaccard similarity to measure how strong the evidence is.

id	claim	evidence
1	(<i>Bill Clinton</i> , <i>President Clinton</i>)	.33
2	(<i>George W. Bush</i> , <i>Dubya</i>)	0
3	(<i>George W. Bush</i> , <i>George H. W. Bush</i>)	.75

Table 2: String similarity as positive evidence

Although string similarity seems to work well in some general cases, its limitations are obvious: i) For claim 2, the evidence based on string similarity has 0 weight, yet we know *George W. Bush* is also known as *Dubya*; ii) For claim 3, *George H. W. Bush* and *George W. Bush* have strong string similarity, yet they are father and son, not a single person.

In our work, we obtain positive evidence to support the claim that *George W. Bush* and *Dubya* refer to the same person. To do this, we need to go beyond string similarity, and explore the relationships among the two instances directly in external sources. Specifically, we explore Wikipedia and the web for positive evidence. Wikipedia has been used for disambiguation in several recent works [13, 9]. In our work, we consider some special constructs used in Wikipedia.

- **Wikipedia Redirects:** Some Wikipedia pages do not have their own content, and accesses to such pages are redirected to other pages. We use $x_i \sim y_i$ to denote the redirection.
- **Wikipedia Internal Links:** Links to internal pages are expressed in shorthand by `[[Title | Surface Name]]` in Wikipedia, where *Surface Name* is the anchor text, and the page it links to is titled *Title*. Again, we denote it as $x_i \sim y_i$, where x_i is the anchor text, and y_i is the title.
- **Wikipedia Disambiguation Pages:** An ambiguous phrase may correspond to multiple Wikipedia pages, each representing a specific interpretation of the phrase. Wikipedia puts such pages together for each ambiguous phrase. We denote this as $x \sim y_i$, where x is the ambiguous phrase, and y_i is the title of any of the Wikipedia pages.
- **Besides Wikipedia,** we also explore positive evidence on the web. We select several patterns, including ‘*x also known as y*’, ‘*x, whose nickname is y*’, etc., and construct a large set of $x \sim y$.

3.3 Quantifying positive evidence

The evidence that support a claim (x, y) may come from multiple sources, and each source gives a score in the range of $[0, 1]$ as an indicator of the strength of the evidence.

Each source employs its own mechanism to score the evidence. For example, for string similarity, a useful measure is the Jaccard coefficient: $w_{(x,y)} = \frac{|x \cap y|}{|x \cup y|}$, where $x \cap y$ denote the set of common words in x and y , and $x \cup y$ denote the union of words in x and y . Alternatively, we can use Dice’s coefficient, $w_{(x,y)} = \frac{2n_t}{n_x + n_y}$, where n_t is the number of character bigrams found in both strings, n_x and n_y are the number of bigrams in string x and y respectively.

Both the Jaccard and the Dice’s coefficients may encounter some problems. Consider the *highschool* category that contain the following instances {*riverdale high school*, *riverfall high school*}. The two have high similarity (0.5 based on Jaccard, 0.9 based on Dice’s coefficient) because they share a substring “high school.” Unfortunately, this is a very common substring in the *highschool* category. To correct this problem, we use weighted Jaccard similarity, where the weight is defined by the inverse of term frequency.

Other sources may collect multiple pieces of evidence for a claim. For example, in Wikipedia links, there are 32,467 occurrences of *united states ~ usa*, and 122 occurrences of *George W. Bush ~ George Bush*. We can use the Sigmoid function to take into consideration the multiple occurrences.

$$w_{(x,y)} = \begin{cases} \frac{1}{1+e^{1-t}} & t > 0 \\ 0 & t \leq 0 \end{cases} \quad (2)$$

where t is the number of occurrences. Hence, if there is only a single piece of evidence backing up a claim, then the evidence has a weight of 0.5. When t becomes larger, the weight becomes closer to 1.

Finally, assume evidence from k sources return a vector of k scores, (a_1, a_2, \dots, a_k) , where each score a_i is in the range of $[0, 1]$. Our goal is to map the k scores into one score in the range $[0, 1]$. We apply the noisy-or model.

$$1 - (1 - a_1) \cdot (1 - a_2) \cdots (1 - a_k) \quad (3)$$

Intuitively, the evidence for a claim has a high score as long as one evidence source gives it a high score. For instance, though *George W. Bush* and *Dubya* share no lexical similarity, Wikipedia frequently suggests *Dubya* is his nickname (strong evidential similarity).

4. METHODS

This section discusses how we use the evidence for entity resolution.

4.1 Problem Definition

We formally abstract claim (x, y) as graph connectivity, using graph $G = (V, E)$ with a set of vertices V representing the union of entities from two taxonomies, and a set of edges E with weight quantifying positive evidence between two entities (Eq 3).

Given G , our aim is to group entities into clusters, so that entities in the same cluster refer to the same real-life entity. Specifically, we want to find a cut $C \subseteq E$ to specify which edges we want to cut from G , such that each connected component in $G'(V, E - C)$ corresponds to the same real-life entity. Among all possible cuts, our goal is to find C that are (a) best supported by positive evidence and (b) not rejected by negative evidence. More formally, our goal is to find G' satisfying the following criteria:

DEFINITION 4.1 (POSITIVE EVIDENCE). $G'(V, E - C)$ should maximize the positive evidence supporting the claim $e \in E - C$ quantified as $w(e)$. In other words, C should minimize the following objective function:

$$\sum_{e \in C} w(e)$$

DEFINITION 4.2 (NEGATIVE EVIDENCE). A valid solution G' should disconnect x and y rejected by some piece of negative evidence N_i .

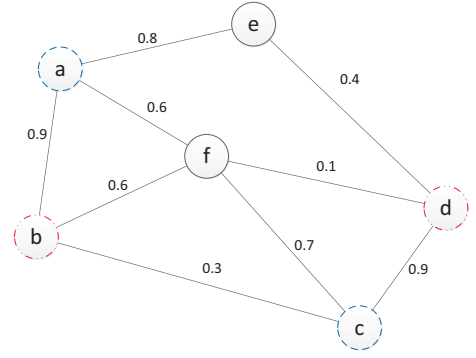


Figure 1: A graph with i) 6 entities connected by positive evidence of different weights, and ii) two pieces of negative evidence: $\{a, c\} \in N_1$ and $\{b, d\} \in N_2$.

This problem can be formalized as the *multi-multiway cut* problem [2]. In this problem, given G and k sets N_1, \dots, N_k of vertices, the goal is to find a subset C of edges whose removal disconnects every $x, y \in N_i$ for some i with minimal cost. This problem can be formulated as the following integer programming problem [2]:

$$\begin{aligned} & \text{minimize} \sum_{e \in C} w(e)x(e) \\ & \text{subject to} \\ & \forall P \in \mathcal{P} : \sum_{e \in P} x(e) \geq 1 \\ & \forall e \in E : x(e) \in \{0, 1\} \end{aligned}$$

where \mathcal{P} denotes the set of all paths between $u, v \in N_i$ for some i . However, finding an exact solution requires ILP (integer linear programming) optimization populating binary decision variables $x(e)$ for every entity pair, $O(10^{14})$ in our problem, with the constraints enumerating all possible paths between every entity pair mentioned in negative evidence and requiring their disconnection. Even the best known approximation algorithm [2] requires expensive LP optimization and \mathcal{P} enumeration, and is still too expensive for our target problem. In addition, this solution has $O(\log k)$ approximation ratio, which is also not desirable for our target problem of large k .

To build a scalable approximation algorithm, we revisit the two principles discussed in Section 3, which we will describe in detail in the following two sections.

- Using the ‘Clean data has no duplicates (**CnD**)’ principle: We use negative evidence from clean data to develop an approximate graph cut that is highly efficient and accurate.
- Using the ‘Birds of a Feather (**BoF**)’ principle: We localize our graph to compare only the entities in *matching categories* c_1 and c_2 with high $\text{sim}(c_1, c_2) > t$, as the **BoF** principle rejects (x, y) , when $x \in c_1, y \in c_2, \text{sim}(c_1, c_2) \leq t$.

4.2 Using the CnD Principle

As Section 3 discussed, negative evidence from some sources, e.g., Freebase created and maintained manually, does not contain any duplicates. Such cleanliness can be exploited for efficient graph cut computation. More formally, we define clean data as follows:

DEFINITION 4.3 (CLEAN EVIDENCE). k sets of negative evidence N_1, \dots, N_k is clean, if and only if $x \in N_i$ and $y \in N_j$ with different names cannot refer to the same real-life entity.

For example, $N'_1 = \{\text{Bill Gates}, \text{Bill Clinton}\}$ and $N'_2 = \{\text{William Gates}, \text{Bill Clinton}\}$ are not clean, as Bill and William Gates refer to the same entity. In clear contrast, $N_1 = \{\text{Bill Gates}, \text{Bill Clinton}\}$ and $N_2 = \{\text{Barack Obama}, \text{Bill Clinton}\}$ is a clean set.

For a clean set, our optimization problem can be reduced to a special case of multi-multiway cut with $k = 1$, known as the *multiway cut* problem [14]. That is, we can aggregate k sets of negative evidence into a single set $\cup_i N_i$. For example, in the above clean set example, we can aggregate N_1 and N_2 into a single set $\cup_i N_i = \{\text{Bill Gates}, \text{Bill Clinton}, \text{Barack Obama}\}$, suggesting any claim (x, y) for $x, y \in \cup_i N_i$ should be rejected. The same aggregation would not work for the dirty set example, as Bill and William Gates in $\cup_i N'_i = \{\text{'Bill Gates'}, \text{'Bill Clinton'}, \text{'William Gates'}\}$ refer to the same person and should not be rejected.

The multiway cut problem is proved to be NP-hard when $m = |\cup_i N_i| \geq 3$, and an efficient “isolation heuristic” gives an approximation ratio of $2 - \frac{2}{m}$ [14]. In this paper, we discuss this heuristic, though there is an approximation algorithm with possibly better ratios, e.g., 1.348 [25], as all other existing solutions do not scale for our graph, requiring LP optimization.

The isolation heuristic works as follows. For any “terminal” $t_j \in \cup_i N_i$, we will find its minimum weight “isolating cut”, which is a subset of edges, whose removal disconnects t_j from all other terminals. Such a cut can be found by connecting all other terminals to a shared new vertex v with infinite weight and running the maximum flow algorithm from t_j to v . Once we find the isolating cut for each terminal node, we can union the sets to get an approximate answer. A more detailed description is available in Algorithm 1.

Algorithm 1 Isolation heuristics

Require: $G = (V, E)$, terminals $\cup_i N_i = \{t_1, \dots, t_m\}$
 For each t_i , compute a minimum weight isolating cut E_i
 Output $E_1 \cup \dots \cup E_{m-1}$, assuming $w(E_1) \geq \dots \geq w(E_{m-1})$
 (without loss of generality)

In our problem, as we collect negative evidence that are both clean and dirty, we divide them into two groups of sets N and N' , i.e., $N \cup N' = \cup_i N_i$ and take a two-phase approach: First, we run the isolation heuristics for clean set N . Once we get connected components, we investigate each $N_i \in N'$, to identify violated entity pairs $s, t \in N_i$ which belong to the same connected component. Once we collect all such pairs, we can find a minimum weight cut using s as a source node and t as a sink, known as $s-t$ cut [18], separating two pairs and thus eliminating a violation. $S-t$ cut can be computed in polynomial time for every violated pair.

For faster computation, we can consider a greedy approximation of isolation heuristics. Specifically, we start from G with singleton clusters, i.e., $G' = (V, \phi)$, where all terminals are trivially isolated. We can then greedily add an isolating cut with the highest weight first.

As the problem is NP-hard, this greedy heuristic finds a sub-optimal solution, as illustrated in Fig 1 with two pieces of negative evidence $N_1 = \{a, c\}$ and $N_2 = \{b, d\}$. Starting from singleton clusters, we iteratively add edges in decreasing order of weights—We will choose to insert edges (a, b) , (c, d) , (a, e) , (c, f) . After these four inserts, the next candidate would be (a, f) , which is not an isolating cut by violating N_1 and generating a path from a to c . Similarly, (b, f) is not an isolating cut by generating a path from

b to d . We thus skip these two candidates and continue to add (d, f) . This insert terminates the search, as we cannot add any more edges without violating negative evidence. This solution with a cost of $0.4 + 0.3 + 0.6 + 0.6 = 1.9$ is sub-optimal, as adding (a, f) and (b, f) instead of (c, f) and (d, f) would lower the cost to $0.3 + 0.3 + 0.1 + 0.7 = 1.4$.

We thus implement a Monte-carlo approach in Algorithm 4.2 to randomize edge insertion with probability proportional to the cost from the “randomized candidate list” *RCL*. Furthermore, this procedure can be repeated several times and the lowest cost cut can be identified as a solution.

Algorithm 2 Monte-carlo heuristics

Require: $G = (V, E)$, $N = N_1, \dots, N_M$
 $E' = \phi$, $RCL = E$
while RCL is not empty **do**
 Randomly select e from RCL with probability proportional to its weight
 $RCL = RCL \setminus \{e\}$
 if e is an isolating cut **then**
 $E' = E' \cup \{e\}$
 end if
end while
 return (V, E')

After a graph cut, each connected component corresponds to one real-life entity. As each Freebase entity corresponds to a unique real-life entity, each component has at most one Freebase entity, which we label with belief 1.0 and then propagate to unlabeled Probase entities using Random Walk with Restart to calculate the label probability of each node in the subgraph. Once the propagation converges, the scores can be used to prune out Probase entities with low scores.

4.3 Using the BoF Principle

As computing cuts is expensive, we reduce the given graph such that entities in E belong to the matching categories. That is, for our aim of entity resolution, a book entity cannot refer to a person. We can thus use ontological information, representing which category the entity belongs to, in order to significantly reduce the graph size.

To decide whether two categories c_1, c_2 are a match, we discussed $\text{sim}(c_1, c_2)$ in Eq. 1, combining the similarity between element sets E_c and attribute set A_c of the two categories. More specifically:

- $f(E_{c_1}, E_{c_2})$: Most existing ontology integration work quantifies a set similarity between E_{c_1} and E_{c_2} , e.g., using Jaccard similarity or its variants as used in [31].
- $f(A_{c_1}, A_{c_2})$: Alternatively, one can compare A_{c_1} and A_{c_2} . From Freebase, each category is associated with a single relational table, from which we can obtain a set of attributes. However, for Probase, we can collect many table instances describing entities of the given category. We thus use *vector space model* used for representing text documents, to represent each category as a frequency vector of a universe of attribute names used, normalized by the number of table instances. For instance, a Probase class *album* is frequently represented by attributes $\{\text{genre}, \text{artist}, \text{producer}\}$, represented by attribute frequency vector $\{0.9253, 0.8431, 0.8301\}$. Meanwhile, attribute frequency vector for the Freebase category will have binary values. The similarity $f(A_{c_1}, A_{c_2})$ is generally computed using cosine similarity or KL-divergence.

Using A_c complements E_c in the following two aspects: *First*, when using E_c similarity alone, two entities with the same name cannot be distinguished, such as *electronics*, which is both an industry and a genre. When comparing two categories of small size, such a *false match* may lead to the overestimation of the concept similarity—In our dataset, two unrelated categories */broadcast/genre* and *manufacturing companies* of size 306 and 108 have a few false matches, such as *electronics* and *automotive*, which generates a high Jaccard similarity score 0.0147. Meanwhile, if considering A_c similarity as well, we can distinguish *electronics* the industry and the genre, as attributes describing the two would be different. *Second*, the size of E_c varies by orders of magnitude over categories and sources, e.g., two identical categories *albums* (from Probase) */music/album* (from Freebase) have 1900 and 526,038 entities respectively. A severe imbalance in the size of E_c , though we attempt normalization, negatively affects the reliability of metrics, while A_c sizes are more balanced.

Accurate computation of both metrics requires to resolve entities and attributes, as treating the entities *Bill Gates* and *William Gates*, or the attributes *author* and *writer*, as unmatched would underestimate the score. However, as we are computing these metrics for the goal of entity resolution, it is unrealistic to assume entities and attributes are resolved *a priori*. We thus consider only exact matches for both entity and attribute similarity, as a lower bounding estimate, to reduce the graph size. Once we compute graph cuts, we can apply our finding, e.g., *Bill Gates* \sim *William Gates*, to recompute $f(E_{c1}, E_{c2})$ so as to refine the score into a tighter lower bound. From this refinement, we can identify some unrelated categories, which were identified as matching, and we can drop entity resolution results obtained from such category pairs. We leave detailed discussions on the details of how we identify similarity threshold t and compute matching pairs efficiently in Appendix C.

5. EXPERIMENTS

In this section, we evaluate our approach for entity resolution. Please refer to Appendix D for a description of the system, including its flowchart, and the setting of the experiments.

5.1 Evidence

We extract positive and negative evidence from Wikipedia and Freebase. Table 3 and Table 4 show the numbers of different types of evidence and their sources. Note that we also use Freebase as a source of negative evidence, as Freebase is handcrafted, so it satisfies the ‘Clean data has no duplicates’ principle.

Wiki Source	# of Pairs	Source	# of Bags
Links	12,662,226	Wikipedia List	122,615
Redirect	4,260,412	Wikipedia Table	102,731
Disambiguation	223,171	Freebase	12,719

Table 3: Positive Evidence

Table 4: Negative Evidence

5.2 Instance Mapping

We match Freebase and Probase categories based on similarity in their attributes and elements. We find 763,000 matchings. Here, we evaluate some selected pairs, including: i) Probase *politician* and Freebase */government/politician*; ii) Probase *format* and Freebase */computer/file_format*; iii) Probase *system* and Freebase */computer/operating_system*; iv) Probase *airline* and Freebase */aviation/airline*. More results can be found in the Appendix (Fig. 5).

We measure precision, recall and output size for each case. Precision and recall are well defined. However, in large-scale instance

mapping, sometimes it is not easy to obtain the exact recall. We use *output size* [1] as an additional factor in our evaluation. When exact recall is not available, the output size can also give us a feeling of how effective the method is in finding all qualified matches.

We manually label instance pairs as ‘Match’ or ‘Non-match.’ We cannot simply use extracted positive evidence for automatic precision evaluation, because, first, it may contain noise; second, it may not contain all variations of possible mappings. Specifically, for a matching category pair, we randomly sample some Freebase entities, then label all Probase entities that map to the Freebase entities as ‘Match,’ and the rest as ‘Non-match.’ Currently, our algorithm maps each Probase entity to one Freebase entity. In our evaluation, for ambiguous instances such as ‘Bush,’ we consider both mappings to ‘George W. Bush’ and ‘George H. W. Bush’ correct.

Table 5 shows the result for the 4 selected category pairs. We labeled 214, 134, 76, 201 Freebase entities for each case, and obtained 1,687 positive mapping pairs (408, 406, 384 and 489 for each category pair, respectively). For our method, we used thresholds 0, 0.05 and 0.10. We compare our method with two baseline algorithms. Let (X, Y) be a category pair. Baseline #1 maps $x \in X$ to $y \in Y$ if (x, y) has the strongest positive evidence (e.g., it has the largest number of $x \sim y$ instances). Baseline #2 maps $x \in X$ to $y \in Y$ if (x, y) has the biggest string similarity as measured by Jaro-Winkler [34], and the similarity is larger than a threshold of .9, as smaller thresholds produce too many false positives.

For popular and stable categories such as *politician* or *airline company*, we have more high quality positive evidence. Therefore, Baseline #1 shows high precision for the pair ‘politician’ and ‘/government/politician,’ because it is based on strong positive evidence. However, Baseline #1 does not take string similarity into consideration. As there are misspelled instances, and variations in surface forms, the recall and the output size of Baseline #1 are smaller than ours. Although Baseline #2 is good for matching names with misspellings, it gives relatively low score to pairs such as (‘Barack Obama’, ‘Obama’), and produces many false positives such as (‘George H. W. Bush’, ‘George W. Bush’), (‘Hillary Clinton’, ‘Bill Clinton’). As a high threshold of .9 is used in Baseline #2 to boost precision, both the recall and the output size of Baseline #2 are smaller than ours, or even Baseline #1’s. On the other hand, the recall of our methods is higher than the baseline methods while the precision is still competitive.

For relatively new and less well-defined categories such as *file format* or *operating system*, using only positive evidence extracted from Wikipedia is not enough for entity resolution. Therefore, both the precision and recall of Baseline #1 is low. Meanwhile, Baseline #2 produced large amount of false positives because the string length is short (for *file format*), or the discriminative part in the string is small (e.g., ‘Windows 95’ and ‘Windows 7’). Our method outperforms both Baseline #1 and #2 in precision, recall, and output size.

In Table 9 (see Appendix), we show a list of Probase entities that are mapped to Freebase instances such as ‘Barack Obama,’ ‘American Airlines,’ and ‘XLS.’ Take the mapping between ‘us president barack obama’ and ‘barack obama’ for example. Baseline #1 failed because it lacks positives evidence from Wikipedia saying ‘us president barack obama’ is ‘barack obama,’ and Baseline #2 failed because the similarity between the two is less than the threshold of .9. Our method worked because the positive evidence from string similarity was not disrupted by any negative evidence. As another example, for ‘American Airlines,’ Wikipedia Disambiguation page provides a wrong piece of positive evidence: ‘Hawaiian airlines’ \sim ‘American Airlines.’ Baseline #1 failed because of this, but our method was able to overcome the noisy evidence with negative ev-

Probase Class Name Freebase Type ID	politicians /government/politician			formats /computer/file_format			systems /computer/operating_system			airline /aviation/airline		
Method	P.	R.	S.	P.	R.	S.	P.	R.	S.	P.	R.	S.
Baseline #1	0.9952	0.6603	603	0.9024	0.2517	1160	0.9040	0.3767	1110	0.9231	0.6326	712
Baseline #2 (0.90)	0.9792	0.3110	451	0.7922	0.2743	1209	0.5965	0.2696	1420	0.9623	0.4788	743
Our method (0.00)	0.9815	0.8413	685	0.9370	0.8605	2222	0.9180	0.5967	5913	1.0000	0.7087	787
Our method (0.05)	0.9851	0.8413	684	0.9654	0.7585	1766	0.9415	0.5367	1981	1.0000	0.7087	782
Our method (0.10)	0.9924	0.8254	677	0.9951	0.6837	1581	0.9568	0.4433	1607	1.0000	0.7087	770

Table 5: Precision and Recall for Selected Category Pairs. (P.: Precision, R.: Recall, S.: Output Size)

idence found in Wikipedia List (using the ‘Clean data has no duplicates’ principle). As yet another example, we have strong positive evidence that ‘Microsoft Excel File’ is ‘XLS,’ weak positive evidence that ‘Microsoft .wav File’ is ‘Microsoft Excel File,’ and through transitivity even weaker positive evidence that ‘Microsoft .wav File’ is ‘XLS.’ Without negative evidence, we might draw the wrong conclusion. However, since the positive evidence is very weak (because the matches are on words Microsoft and File, which have high frequency and hence low score), such mistakes can be filtered out by a threshold as low as 0.005.

5.3 Finding Candidate Category Pairs

We demonstrate how Eq 1 implements the BoF principle in finding candidate category pairs. First, Table 6 shows that Jaccard similarity alone does not work well for two reasons: i) The two sets have a big difference in size; and ii) Instances have ambiguous names. More specifically,

- A large Jaccard similarity threshold will exclude related pairs such as (Written Work, Novels) because although the categories are heavily related, they may not have enough similarity due to a large size difference.
- A small Jaccard similarity threshold will fail to exclude false positives such as (Places, Written Work). This is because there are instances of Written Work, Musical Album titled China, Canada, etc.

Second, as the first 4 rows in Table 6 show, attribute similarity fares better than entity similarity in distinguishing ambiguous category pairs, allowing us to obtain (Novels, Written Work) and (Albums, Musical Albums), and reject (Places, Written Work) and (Words, Musical Album). However, since Freebase has only a few attributes, using attributes similarity alone may also cause problems.

Overall, a linear combination of entity similarity and attribute similarity works well. Besides using similarity measures only, we can also use the hierarchy of Probase and a few handcrafted rules to improve the precision and recall in matching.

5.4 Scalability

In our work, we integrate positive and negative evidence by finding a multiway cut in a connected graph. A common approach to find a multiway cut is to use integer linear programming (ILP). To estimate the cost of doing this for integrating web scale taxonomies, we select three categories (of typical size) from Probase, connect entities in each category by positive evidence, and find the largest connected subgraphs. Table 8 shows that typically the largest connected subgraph contains almost half of the nodes in the original graph. Since we need to apply ILP on each connected subgraph, the sheer size of the subgraph makes the ILP approach computationally infeasible [2]. This justifies the use of a Monte Carlo approach to solve our problem.

Probase category	$ V_1 $	$ V $	$ V_1 / V $
song	9932	21963	0.4522
artist	18839	39648	0.4752
disease	1435	4480	0.3203

Table 8: Size/ratio of the largest connected subgraph. V_1 is the set of nodes in the largest connected subgraph of $G = (V, E)$

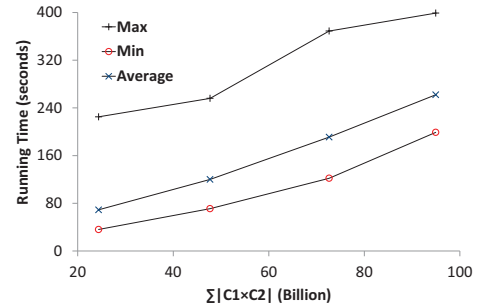


Figure 2: Running Time

Next, we show the performance of our method. For the input, we select four sets of category pairs. For each set, the total size of Cartesian product between each pair is about 23 billion, i.e., $\sum |C_1 \times C_2| \approx 23$ billion, where (C_1, C_2) is a pair of categories (Table 7).

We measure the running time for input set 1, the union of input set 1 and input set 2, and so on. As we process instance mapping in parallel, we measure the running time of each process and then compute the maximum, the minimum, and the average of them. Fig. 2 shows the result. We see that the average running time is about 2 minutes to process 48 billion instance pairs. The reason why the maximum running time of an instance mapping client is much larger than the average is that the size of Freebase types has a very biased distribution: The Book (/book/book), Person (/people/person), and Location (/location/location) categories contains 2.4 million, 1.5 million, and 1.2 million instances respectively whereas most other types contain much fewer instances. However, as the number of pairs increases, the experiment shows that our method demonstrates almost linear performance.

6. CONCLUSION

Entity resolution for data integration is a challenging task. In this paper, we study the problem of matching millions of entities in two web scale taxonomies. Unlike integrating two relational tables, taxonomies may not contain much information about each entity. But it is exactly this reason that makes the task of integrating two taxonomies important, as integration serves as an indispensable mechanism for taxonomies to enrich themselves by ‘‘borrow-

Freebase Type	Probase Class	$ E_{c_1} \cap E_{c_2} $	$f(E_{c_1}, E_{c_2})$	$f(A_{c_1}, A_{c_2})$	$sim(c_1, c_2)$
Written Work	Novels	755	0.0004	0.0275	0.0220
	Places	4902	0.0023	0.0002	0.0006
Musical Album	Albums	1095	0.0026	0.0638	0.0516
	Words	1550	0.0036	0.0007	0.0013
Breed Origin	Country	82	0.0826	0.0000	0.0165
Musical Instrument	Percussion	95	0.0565	0.0035	0.0141

Table 6: Similarities of Class Pairs. $\lambda = 0.2$

	$\sum C_1 \times C_2 $
Task Set 1	24,377,292,299
Task Set 2	23,323,215,795
Task Set 3	24,923,006,886
Task Set 4	22,313,005,928

Table 7: Size of Each Task Set

ing” content from other taxonomies. We develop a framework that relies on the interconnections of the data in the taxonomies as well as in external data sources for entity resolution. We collect a large number of positive and negative evidence from the interconnections, and formulate the task of entity resolution as a multi-way graph cut problem. Our experiments show that our method scales up to millions of categories and entities, and produces very high quality resolutions.

7. REFERENCES

- [1] A. Arasu, M. Götz, and R. Kaushik. On active learning of record matching packages. In *SIGMOD*, 2010.
- [2] A. Avidor and M. Langberg. The multi-multiway cut problem. In *SWAT*, pages 273–284, 2004.
- [3] R. J. Bayardo, Y. Ma, and R. Srikant. Scaling up all pairs similarity search. In *WWW*, pages 131–140, 2007.
- [4] I. Bhattacharya and L. Getoor. Collective entity resolution in relational data. *TKDD*, 1:2007, 2006.
- [5] I. Bhattacharya and L. Getoor. Entity resolution in graphs. *Mining graph data*, page 311, 2006.
- [6] I. Bhattacharya and L. Getoor. A latent dirichlet model for unsupervised entity resolution. In *SIAM Data Mining*, 2006.
- [7] M. Bilenko and R. J. Mooney. Adaptive duplicate detection using learnable string similarity measures. In *SIGKDD*, 03.
- [8] K. Bollacker, C. Evans, P. Paritosh, T. Sturge, and J. Taylor. Freebase: a collaboratively created graph database for structuring human knowledge. In *SIGMOD*, 2008.
- [9] R. Bunesco. Using encyclopedic knowledge for named entity disambiguation. In *In EACL*, pages 9–16, 2006.
- [10] S. Chaudhuri, V. Ganti, and R. Kaushik. A primitive operator for similarity joins in data cleaning. In *ICDE*, 2006.
- [11] P. Christen. A comparison of personal name matching: Techniques and practical issues. *Data Mining Workshops, International Conference on*, 0:290–294, 2006.
- [12] W. W. Cohen, P. Ravikumar, and S. E. Fienberg. A comparison of string distance metrics for name-matching tasks. In *IJCAI*, pages 73–78, 2003.
- [13] S. Cucerzan. Large-scale named entity disambiguation based on wikipedia data. In *In Proc. 2007 Joint Conference on EMNLP and CNLL*, pages 708–716, 2007.
- [14] E. Dahlhaus, D. S. Johnson, C. H. Papadimitriou, P. D. Seymour, and M. Yannakakis. The complexity of multiterminal cuts. *SIAM Journal on Computing*, 23:864–894, 1994.
- [15] A. Doan, Y. Lu, Y. Lee, and J. Han. Object matching for information integration: A profiler-based approach. In *Proc. of IIWEB*, pages 53–58, 2003.
- [16] A. Doan, Y. Lu, Y. Lee, and J. Han. Profile-based object matching for information integration. *IEEE Intelligent Systems*, pages 54–59, 2003.
- [17] X. Dong, A. Halevy, and J. Madhavan. Reference reconciliation in complex information spaces. In *SIGMOD*, pages 85–96. ACM, 2005.
- [18] J. Edmonds and R. M. Karp. Theoretical improvements in algorithmic efficiency for network flow problems. *J. ACM*, 19:248–264, April 1972.
- [19] A. Elmagarmid, P. Ipeirotis, and V. Verykios. Duplicate record detection: A survey. *TKDE*, pages 1–16, 2007.
- [20] C. Galvez and F. Moya-Anegón. Approximate personal name-matching through finite-state graphs. *Journal of the American Society for Information Science and Technology*, 58(13):1960–1976, 2007.
- [21] J. Gong, L. Wang, and D. W. Oard. Matching person names through name transformation. In *CIKM*, 2009.
- [22] Google. Freebase data dumps. <http://download.freebase.com/datadumps/>, 2010.
- [23] X. Han and J. Zhao. Named entity disambiguation by leveraging wikipedia semantic knowledge. In *CIKM*, 2009.
- [24] M. A. Hearst. Automatic acquisition of hyponyms from large text corpora. In *COLING*, pages 539–545, 1992.
- [25] D. R. Karger, P. Klein, C. Stein, M. Thorup, and N. E. Young. Rounding algorithms for a geometric embedding of minimum multiway cut. *STOC ’99*, pages 668–678.
- [26] A. Monge and C. Elkan. The field matching problem: Algorithms and applications. In *SIGKDD*, 1996.
- [27] S. Sarawagi and A. Bhamidipaty. Interactive deduplication using active learning. In *SIGKDD*, pages 269–278, 2002.
- [28] P. Singla and P. Domingos. Entity resolution with markov logic. In *ICDM*, pages 572–582, 2006.
- [29] S. Tata and J. M. Patel. Estimating the selectivity of tf-idf based cosine similarity predicates. *SIGMOD Rec.*, 36:75–80, December 2007.
- [30] S. Tejada and C. A. Knoblock. Learning domain-independent string transformation weights for high accuracy object identification. In *SIGKDD*, pages 350–359, 2002.
- [31] O. Udrea. Leveraging data and structure in ontology integration. In *SIGMOD*, pages 449–460. ACM Press, 2007.
- [32] Y. Wang, H. Li, H. Wang, and K. Q. Zhu. Toward topic search on the web. Technical report, Microsoft Research, 2010.
- [33] S. Whang, O. Benjelloun, and H. Garcia-Molina. Generic entity resolution with negative rules. *The VLDB Journal*, 18(6):1261–1277, 2009.
- [34] W. E. Winkler. The state of record linkage and current research problems. Technical report, Statistical Research Division, U.S. Census Bureau, 1999.
- [35] W. Wu, H. Li, H. Wang, and Q. Zhu. Towards a universal taxonomy of many concepts. 2010.
- [36] C. Xiao, W. W. 0011, X. Lin, and J. X. Yu. Efficient similarity joins for near duplicate detection. In *WWW*, pages 131–140, 2008.

APPENDIX

A. RELATED WORK

Entity resolution, also known as record linkage or reference reconciliation, is an important and difficult problem. Some researchers try to solve this problem by finding the best string similarity measures [26, 12]. But string similarity measures have limitations: They are not able to identify nicknames, or distinguish entities whose names are very similar, or handle noises in names. Recently, researchers focused on adaptive algorithms that learn similarity measures automatically [7, 30]. However, the difficulty of using these methods cannot be overlooked: It requires training data for each domain. To address this problem, some approaches use active learning so that only a little user input is needed [27]. But, still no methods can effectively handle the case where an entity can have totally different names, for example, ‘The Governor’ as a nickname for ‘Arnold Schwarzenegger.’ Also, labeling a small number of examples for each domain is still costly in our scenario, since the taxonomies we are dealing with contain millions of categories (domains).

Entity resolution in specific domains is challenging as well. To match person names [20, 11, 21], a variety of rules, such as abbreviation, omission, transposition (or sequence changing), are developed to transform a person’s name before comparing it with other names [21]. Profilers that have domain knowledge (for domains such as movies, reviews, and people) are used for various object matching tasks [16, 15]. More recently, Markov logic is used [28] to soften the hard constraints of handcrafted rules in first-order logic. Because these methods depend on much domain knowledge, they work well for specific domains only.

Entity resolution is also studied extensively in the field of data integration [19]. A database table usually has multiple columns and contains many values. In this setting, some attributes may be more important than others for measuring similarity for records. Arasu et al [1] proposed a method that learns threshold-based boolean functions or linear classifiers for entity resolution. Given two tables and a precision threshold, their method provides the largest output (a concept similar to recall defined in their paper) satisfying the precision threshold. There are many challenges. First, it is difficult to construct a good training set that includes both positive and negative training data. Second, the training sets are usually domain-dependent. Third, attributes in the two tables may require reconciliation as well. A possible approach is to regard attributes and their values as another type of classes and entities, and leverage relations among them to help solve the entity resolution problem [5, 6, 4].

Recent approaches also leverage knowledge acquired from external sources for domain-independent entity resolution. For example, some extracts entities’ surface forms (names or aliases) from Wikipedia and builds a synonym dictionary for entities [13, 9, 23]. Given a name, we first find Wikipedia articles corresponding to the name or to its synonyms, and then create a bag-of-words vector for the name, using the words from the Wikipedia articles. Finally, we compare two names using metrics such as the cosine similarity for entity resolution. This approach assumes entities have corresponding Wikipedia articles, but that covers a very small percentage of entities. One way to extend it is to consider transitive relation of positive evidence (i.e., if x is a synonym of y , and y is a synonym of z , then x might be a synonym of z), although it may introduce some noises.

Some recent work employs *negative evidence* for entity resolution. Dong et al [17] focused on the case where each entity has a set of attribute values. They use a propagation method to link related mentions, and use negative evidence to restrict the propaga-

tion. However, the constraint rules for negative evidences is handcrafted. Most recently, Whang et al [33] use domain knowledge as negative rules. Entity resolution relies on attributes and values of the entities, and the approach does not scale well when the number of entities and rules become large.

B. THE PROBASE TAXONOMY

Fig. 3 shows an interface with which users can browse the Probase taxonomy, and we can see a category (politicians) has many super categories, sub categories, instances, and similar categories.



Figure 3: The Probase Taxonomy.

C. CLASS MATCHING ALGORITHM

For deciding threshold t defining matching E_c or A_c , we model both E_c and A_c as frequency vector and estimate the distribution of *dot product similarity*, as it is the common numerator in both Jaccard similarity and cosine similarity. For each Freebase c_1 , the distribution of its dot product similarity to all Probase classes, according to the observation in [29], can be modeled as inverse normal distribution and its mean and variance can be estimated from the mean and variance obtained from Probase distribution. This distribution is characterized by a mass of probability close to zero, corresponding to the pairs unmatched, with the minority long tail with high score, corresponding to matching pairs. We can identify t with drastic probability change as threshold.

Once t is identified, a naive solution for finding a matching class with the highest cosine similarity score would be all-pair computation, for $12.7 \text{ thousand} \times 2 \text{ million}$ class pairs. Though a naive all-pair class comparison may be feasible for this specific case with a small number of categories in Freebase, it cannot scale for joining two large-scale ontologies, each with millions of categories (or, for self-joining Probase).

For more scalable computation of $f(A_{c_1}, A_{c_2})$, we can build inverted indices that map each attribute to a list of Probase classes (sorted in the order of probability), to efficiently locate the matching class and minimize computation of the similarity score. Such an index enables to locate the matching Probase class with sub-linear scan for each Freebase class and thus significantly reduce pairwise computation.

For $f(E_{c_1}, E_{c_2})$, we can view the problem as a *set similarity join* by the overlap of two sets X and Y . Many efficient algorithms have been proposed [36, 10, 3], which order all sets based on the same ordering. Once ordered, two sets with high enough similarity would share a prefix, and the length of the shared prefix can be used to safely prune out class pairs with short length. Using this principle, these algorithms drastically reduce the candidate pairs and are readily applicable to our problem to achieve scalability.

Freebase Entity	Baseline #1	Baseline #2 (0.90)	Our Method (0.10)
Barack Obama	barack obama, barrack obama, senator barack obama, president barack obama	barack obama, barrack obama	barack obama, barrack obama, senator barack obama, president barack obama, us president barack obama, mr obama
John Kerry	john kerry, senator john kerry, sen. john kerry, senator kerry	john kerry	john kerry, senator john kerry, sen. john kerry, senator kerry, massachusetts sen. john kerry, sen. john kerry of mas- sachusetts
MP3	mp3, mp3s, mp3 files, mp3 format	mp3, mp3s	mp3, mp3s, mp3 files, mp3 format, high-quality mp3, mp3 songs
XLS	xls, microsoft excel	xls, xlsx	xls, microsoft excel, excel file, excel documents
Windows VISTA	windows vista, windows vista sp2, win- dows vista service pack 2	windows vista, win- dows vista sp2	windows vista, windows vista sp2, win- dows vista service pack 2, microsofts windows vista
American Airlines	american airlines, american airline, aa, hawaiian airlines	american airlines, american airline	american airlines, american airline, aa

Table 9: Probase instances mapped to the Freebase entities.

D. EXPERIMENT SETTING

We conduct comprehensive experiments in integrating Probase and Freebase. We build a distributed entity resolution system (as Fig. 4 shows) that contains three servers running 64-bit Microsoft Windows Server 2003 Enterprise SP2 OS, with 16 core 2.53 GHz Intel Xeon E5540 processors and 32 GB of memory. Each of the servers run 10 instance mapping clients. All together, we run 30 instance mapping clients in parallel.

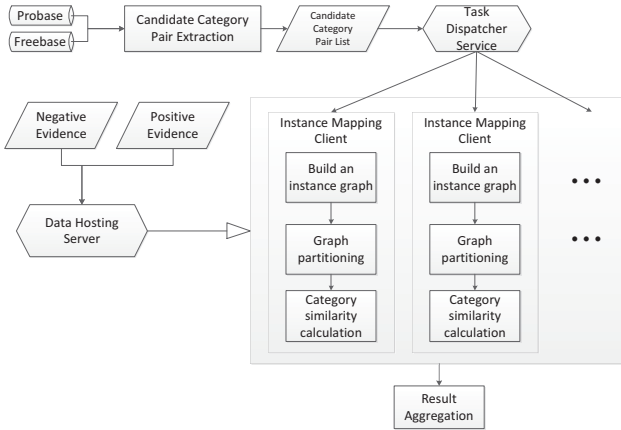


Figure 4: The Flowchart of the System

We used Freebase data dump for 2010-10-14 [22]. The total number of topical instances in Freebase is 13,491,742. Non-topical instances are instances that do not represent real world entities, for example, data types that are used internally in Freebase, and we do not include them in our experiments. The number of non-empty topical categories in Freebase is 12,719.

In terms of Probase, instances are classified into a very diverse set of categories. In the version of Probase we used for the experiments, the total number of instances is 16,423,710, the total number of categories is 2,026,806, and the total number of attributes is 1,727,580.

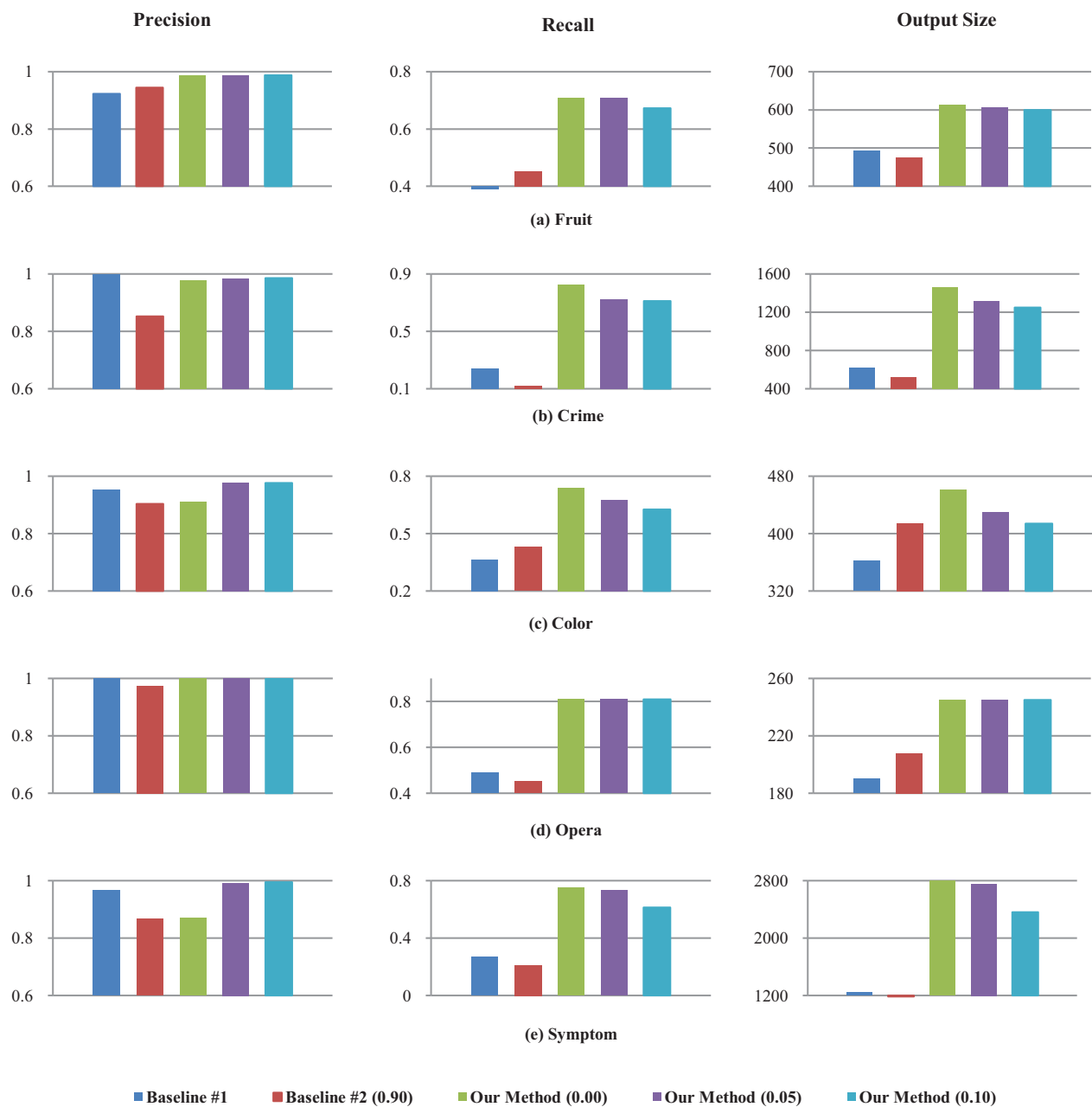


Figure 5: Precision, Recall and Output Size

Category Pair Name	Probase Class Name	Freebase Type ID
Fruit	fruits	/food/ingredient
Crime	crimes	/base/fight/crime_type
Color	colors	/visual_art/color
Opera	operas	/opera/opera
Symptom	symptoms	/medicine/symptom

Table 10: Category Pair Name, and Its Probase Class Name and Freebase Type ID