# Modeling Virtual Footprints

Rajiv Kadaba, Suratna Budalakoti, David DeAngelis, and K. Suzanne Barber

The Laboratory for Intelligent Processes and Systems
The University of Texas at Austin
University Station C5000, ACE 5.124
Austin, Texas, 78712-0321 USA
{kadaba,sbudalakoti,dave,barber}@lips.utexas.edu

**Abstract.** Entities interacting on the web establish their identity by creating virtual personas. This research models identity using the *Entity-Persona Model* which is a semantically annotated social network inferred from the persistent traces of interaction between personas on the web. A *Persona Mapping Algorithm* is proposed which compares the local views of personas in their social network referred to as their *Virtual Signatures*, for structural and semantic similarity. The semantics of the social network of the *Entity-Persona Model* is modeled by a vector space model of the text associated with the personas in the network, which allows efficient comparison of their *Virtual Signatures*. This enables all the publicly accessible personas of an entity to be identified on the scale of the web. This research enables an agent to identify a single entity using multiple personas on different networks. The agent is able to increase the trustworthiness of on-line interactions by establishing the identity of entities operating under multiple personas. Consequently, reputation measures based on on-line interactions with multiple personas can be aggregated and resolved to the true singular identity.

**Keywords:** trust, social networks, identity management, virtual signatures

## 1 Introduction

The way that an individual's identity is created and experienced is fundamentally different in the virtual world. The basic cues used to uniquely identify individuals in the real world are missing, making the association between an entity and its identity ambiguous [14]. This research creates a model of the virtual world which dispels this ambiguity, allowing the virtual personas created by an entity to be linked together. Informally, a virtual persona is a name and its associated attributes, which an entity uses to communicate with other personas.

The virtual world in the context of this research refers collectively to the various explicit or inferred social networks on the web. Examples of explicit social networks are websites such as Facebook, Orkut, MySpace, and LinkedIn. Social networks can be inferred from the digital traces of interaction between entities, or individuals, on the internet, such as in the Blogosphere [9], Online

95

Discussion Forums, Knowledge sharing sites, IRC Logs and the co-occurrence of names in the large amount of textual data on the internet [7]. In an explicit social network there exists a framework by which entities can specify to whom they are related and the context of this relationship. Access to explicit networks is generally controlled [6] because of the privacy concerns of its participants [1].

Inferred social networks lack the privacy mechanisms of explicit networks as its users assume they are as anonymous as they wish to be. The work in this paper counters this assumption since users must engage in information rich interactions in order to provide value to the framework. The establishment of the reputation of an entity's virtual persona within the framework is an important motivating factor for its consistent use, as others use reputation to assess the reliability of information associated with the persona [5]. Every new persona created will need to establish its reputation within its social network which requires time and effort. This penalty associated with creating a persona which is capable of meaningful interaction makes a persona valuable. Virtual personas with erratic interactions are not worth detecting as they have little value.

Search engines treat personal names and pseudonyms as keywords, giving virtual personas the same status as ordinary text. Queries for people's names only find occurrences with verbatim matches to the query text while they may have interacted extensively using various personas. The model proposed in this paper can be used to find more accurate results for the information associated with an individual available publicly on the web. Augmenting web search with the ability to link entities and their personas can be perceived as an attack on an individual's privacy, as information which may have been exchanged with the expectation of anonymity granted by a virtual persona is now linked back to its progenitor. Conversely this research can also contribute to an individual's ability to safeguard their privacy and protect their identity from theft. As the concept of identity in the virtual world is formalized and an upper bound on the ability of a determined adversarial agent is found, techniques to remain anonymous in spite of sophisticated statistical tools can be developed. Further, software agents who are capable of warning users of the unintended inferences which can be made with the data they publish may also be possible.

Another application of this research is in the task of anti-aliasing in social networks. Anti-aliasing [11] is the task of identifying when a single user has multiple aliases in a social network. This is important in trust networks, as a distrusted user or set of users, after being removed from a network (by a moderator, for example), can insert themselves into a network at a later point in time with a pseudonym. Agents capable of continuously associating personas with the singular identity of an entity will offer increased assurance of entity reputations based on interactions of associated personas. Consequently, agents can assess trustworthiness of individual personas and relate those trustworthiness assessments to the singular identity associated with those personas. This research offers the groundwork for establishing the connection between the multiple personas and a singular entity identity.

## 2 Related Work

The privacy of individuals, referenced in social network data released to researchers, application developers, and marketing organizations is ostensibly protected by anonymizing the social graph, as their goal is to make inferences about the aggregated data not specific individuals. Privacy preserving data mining is a research area in which data sets are modified and algorithms developed which do not compromise privacy [16]. Approaches to the reverse task of de-anonymization are usually based on unique subgraphs in the anonymized social network, and are classified into active and passive attacks [3]. Active attacks are attacks in which nodes that form a unique subgraph are inserted into the graph before it is anonymized (for example, before being released to the public for research of other purposes). As this subgraph is known a priori it can be used to identify other nodes in the released graph. Passive attacks [10] are similar, however no nodes are inserted, instead a small group of nodes collude to generate a subgraph which is later used to re-identify adjacent nodes. De-anonymization of individuals in an anonymized data set is equivalent to linking personas in two different social networks as the attackers background information is also a social network. These techniques de-anonymize nodes using structural properties of graphs and are not designed to target specific nodes. They also rely on the fact that both networks have many nodes in common. Purely structural equivalence techniques break down if the neighborhood of an individual node changes drastically. The work presented here uses structural information in addition to persona content for de-anonymization.

In contrast to graph-based approaches, Novak et al. [11] use a content-based approach where they reconcile online personas to unique users by clustering the content associated with the personas, such that each cluster represents a unique user. The data is derived from an online discussion board by only considering text associated with a persona independent of relation information. Similar to the approach used in this paper, two sets of personas are synthetically created by random division. Random division allows words from specific topics to appear in every division which will not occur in real data, making the entity disambiguation results flawed. Algorithms must be robust with respect to text originating from very specific topics which are never repeated. Temporal division as used here addresses this problem.

Jin et al. [7] extract social networks from the web using hypertext data retrieved from search engine queries. The nodes in the network are named entities which are known a priori. Edges are inferred using heuristics from co-occurrence of names, and the type of relationship the edge represents is determined from the query text. Queries consist of the entity names and the type of relationship. Staddon et al. [13] have leveraged web search to determine unintended inferences which can be drawn from data published on the web. Keywords are extracted from data intended to be published using TF-IDF[1] and are used to construct queries to search engines. New keywords in the results of these queries not present

---

[1] Term Frequency - Inverse Document Frequency

in the original keyword set represent inferences which can be drawn from the web. Both these techniques only find information associated with a specific name i.e. they assume an individual has only one persona on the internet. Although the results are interesting they bring little insight to the true nature of identity on the web.

## 3 Modeling Entities and Personas

Entities communicate on the web using identifiers unique within a framework making the identifier - framework combination also unique. At internet scope the identifier can be an IP or email address depending on the protocol, or a username in the scope of a Web 2.0 application. The identity of an entity possessing an identifier is characterized by the set of other identifiers it has communicated with and the content of this communication, collectively referred to as its virtual signature.

### 3.1 Model and Definitions

**Definition 1.** *An entity $\xi$ is something capable of independent interaction, which can be uniquely identified in the real world.*

An entity can be an individual or software agent. Individuals by default are unique as they can have only one instance in the real world. A software agent may not be unique as it is very easy to create many instances of the same agent. Therefore, all instances of the same software agent which exhibit the same behavior are considered collectively a single entity regardless of their physical location. Entities interact on the web through a framework using a persona, which is an instantiation of an entity within the framework. An entity can possess more than one persona within a given framework.

**Definition 2.** *A framework is an implementation of software and associated protocols which enables entities to interact.*

**Definition 3.** *A persona $\pi$ is a tuple $(i, d)$, where $i$ is an identifier unique within a framework and $d$ is a n-tuple of associated information and attributes which an entity uses to establish its identity and interact within a social network. For every persona there exists exactly one entity*

A social network (Definition 4) is used to model the virtual signatures of entities or personas which are nodes in the network such that two nodes are connected if they have communicated. Again, the criteria for considering that communication has occurred depends on the framework.

**Definition 4.** *A social network $S$ is a vertex and edge labeled undirected graph $G = (V, E)$, where $V$ is a set of either exclusively entities or personas and $E$ is a subset of the cartesian product $V \times V$ such that $(v, v) \not\subset E$ [2]. Every edge*

---

[2] No self loops are allowed since they are meaningless in a communication graph.

$(u, v)$, $u, v \in V$ *has a label* $\gamma$ *and every vertex has a label* $\chi$, *which is arbitrary information associated with the edge or vertex.*

The *Entity-Persona Model* consists of a social network of entities $S_\Xi$ and at least one social network of personas $S_\Pi$. $S_\Xi$ is a real world social network of entities whose personas need to be linked. $S_\Pi$ is a social network of personas inferred from the records of the framework in which the entities in $S_\Xi$ communicate. The label of vertex $u_\pi$, $\chi_\pi$ in the graph $G_\Pi$ of $S_\Pi$ is $\pi(i, T(d))$ and the label $\gamma_{\pi_1 \pi_2}$ of edge $(u_{\pi_1}, v_{\pi_1})$ between vertex $u_{\pi_1}$ labeled by $\pi_1$ and vertex $v_{\pi_1}$ labeled by $\pi_2$ is $T(d_{\pi_1} \cup d_{\pi_2})$. $T$ is an operation which reduces $d$ to its semantics which is outlined in section 3.3. Hence, the social network of personas is vertex labeled by the semantics of information generated by the personas and edge labeled by the semantics of all the information exchanged by the personas it connects.

**Definition 5.** *The virtual signature of a persona* $\pi$ *in the context of the* Entity-Persona Model *is the social network* $S_\Pi$ *rooted at the node of the persona.*

The identity of a persona is defined by the structure of $S_\Pi$ and the semantics of the graph labels. $S_\Pi$ is not unique to a persona as it is a social network of personas, however its local view can be unique making its virtual signature also unique.

## 3.2 Problem Specification

Linking entities and their corresponding personas can be formally expressed as finding a mapping between a set of $n$ entities $\Xi = \{\xi_1, \xi_2, ..., \xi_n\}$ which are nodes in a social network $S_\Xi$ and a set of $m$ personas $\Pi = \{\pi_1, \pi_2, ..., \pi_m\}$. $\Pi$ may be partitioned into $P = \{\Pi_1, \Pi_2, ..., \Pi_l\}$ a collection of $l$ subsets of $\Pi$ if the personas exist in $l$ different frameworks. Each framework is expressed as a separate social network in the model.

**Definition 6.** *A mapping* $\mu$ *is a binary relation on* $\bigcup\limits_{i=1}^{l} S_E \times S_{\Pi_i}$, *where* $(\xi, \pi) \in \mu$, $\xi \in S_E$, $\pi \in S_{\Pi_i} \iff$ *entity* $\xi$ *has created persona* $\pi$.

From the information exchanged within each framework, $P$ can be inferred. A query on the model is a set of entities $\Xi$ and a mapping $\mu$ which maps every entity in $\Xi$ to at least one persona in $P$. If no mapping between an entity and a persona is known, the query algorithm will not have an example to train on making mapping the personas of the entity impossible. The social network $S_\Xi$ and the mapping $\mu$ are only partially known. From this information the complete social network $S_\Xi$ and $\mu$ must be inferred. This can be accomplished by,

1. Selecting a entity $\xi \in \Xi$.
2. For $(\xi, \pi_\xi)$ in mapping $\mu$, find a set $\Pi_{new}$ which contains all $\pi \in P$ that best match $\pi_\xi$ by some objective criteria.
3. Update mapping $\mu$ such that $\mu = \mu \bigcup\limits_{q \in \Pi_{new}} (\xi, q)$.
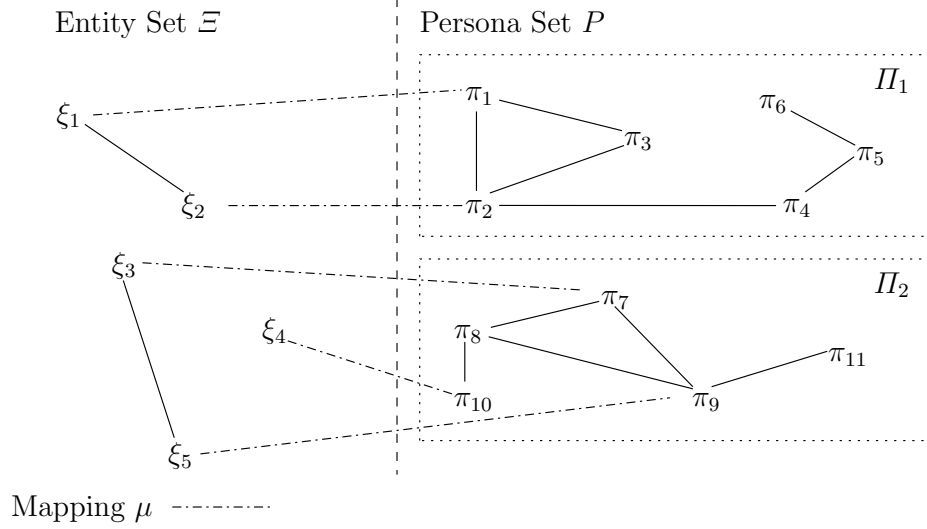
**Fig. 1.** The entity-persona model as it is initilalized.

4. Update $S_\Xi$ such that if two mapped personas have an edge between each other, their corresponding entities have an edge in $S_\Xi$.
5. Repeat.

### 3.3 Semantic Similarity

If the virtual signatures (Definition 5) of personas are to be compared, it follows that there must be a way to compare their local views of their social network. Graphs can be compared structurally using combinatorial algorithms [4] or by comparing their spectra [15]. In a social network however the nodes and edges are labeled requiring the development of the means to compare these labels.

The edge labels of a social network are the inferred semantics of the relationship the edges represent. The vertex labels are the inferred semantics of all the information generated by an entity through the persona associated with the vertex. These semantics must be short descriptions of the information they are inferred from, while maintaining as much discriminative information as possible.

The simplest technique to infer the semantics of text associated with a persona is to construct a vector of unique terms in the text and assign it to a graph element label. Labels can be compared by using the Jaccard index which measures the similarity between two sets $\chi_1$ and $\chi_2$.

$$J(\chi_1, \chi_1) = \frac{|\chi_1 \cap \chi_2|}{|\chi_1 \cup \chi_2|} \tag{1}$$

This produces a score between 0 and 1, with 1 being exactly the same. It is good at using stylometric information such as consistent misspellings to discriminate

between personas but cannot capture information such as favorite words. It is particularly weak at discriminating when the number of terms in the original text is very small.

Term Frequency - Inverse Document frequency is a popular technique in information retrieval to infer which terms best represent a document in a corpus. In the case of the *Entity-Persona Model*, the document is the textual information associated with a persona and the corpus is the collection of all personas. This technique can reduce the dimensionality of the feature space to a greater extent than the previous approach as only the most important terms can be considered i.e. terms with the highest tf-idf. To compare labels $\chi_j$ and $\chi_k$ the cosine similarity between their vectors is computed,

$$\text{Cosine Similarity}_{\chi_j, \chi_k} = \frac{\sum_{\forall i} tf - idf_{\chi_j, t_i} \times tf - idf_{\chi_k, t_i}}{||tf - idf_{\chi_j}|| \times ||tf - idf_{\chi_k}||} \tag{2}$$

Cosine similarity results in a score between $-1$ and $1$ with $1$ being exactly the same. Therefore a higher score implies that the labels being compared are similar.

### 3.4 The Persona Mapping Algorithm

The algorithm takes as input, a source graph $G_S$ of a social network of personas $S_S$, a source vertex $s \in G_S$ of a persona $\pi_s$, and a target graph $G_T$ of a social network of personas $S_T$. $S_S$ and $S_T$ can be the same social network, in this case we are looking for an entity with multiple personas in the same social network. It returns a vertex in target graph which has the maximum overlap between its virtual signature and the signature of $s$. The algorithm is a modified simultaneous *Breadth First Search* of two graphs. It also takes as inputs $\epsilon$ and $\Delta$ which help prune the search space, and $\sigma$ which is the standard deviation of a gaussian used to weigh the importance of a semantic or structural information while producing a score of similarity between two personas.

As the personas of an entity may not have the same neighbors in their social graphs, the structure of the local view of their social network can vary. Depending on the framework the personas are inferred from they may exchange information on very different topics making semantic information unreliable. If both semantic and structural information change drastically a persona may not be able to be mapped correctly, but if only one of them change the algorithm can take this into account by weighing both types of information differently using a zero mean gaussian to decide on the weight of the scores of vertices in the auxiliary graph. The larger the path length between the vertex which maps the source vertex to the candidate vertex and a vertex in the auxiliary graph, the less that vertex contributes to the score. If $s_{i,d}$ is the score of the $i^{th}$ vertex at path length $d$ then the score of the similarity of the candidate vertex to the source vertex is given by:

$$score = \sum_{d=0}^{\Delta} \frac{\sum_{\forall i} s_{i,d}}{|s_{i,d}|} \left( \frac{1}{\sqrt{2\pi}} e^{-\frac{d^2}{2}} \right) \tag{3}$$

---
**Algorithm 1** Persona Mapping Algorithm.
---

MAP-NODE($G_S, G_T, s, \epsilon, \sigma, \Delta$)

```
 1   Q ← ∅
 2   for each vertex v ∈ V[G_T]
 3          do Q ← GET-VERTEX-SCORE(G_S, G_T, s, v)
 4   while Q ≠ ∅
 5          do VertexScore ← MAXIMUM(Q)
 6              c ← EXTRACT-MAX(Q)
 7              Q_S, Q_T ← INITILIZE-QUEUES(s, c)
 8              G_A ← INITILIZE-GRAPH(s, c, VertexScore)
 9              while Q_S ≠ ∅ and CHECK-DEPTH(Δ)
10                  do node_1 ← DEQUEUE(Q_S)
11                      node_2 ← DEQUEUE(Q_T)
12                      AssignedEdges ← MAP-EDGES(node_1, node_2)
13                      UPDATE-GRAPH(G_A, AssignedEdges)
14                      Q_S, Q_T ← UPDATE-QUEUES(Q_S, Q_T, AssignedEdges)
15              Score ← GET-COMBINED-SCORE(G_A, σ)
16              if LastScore − Score ≥ ε
17                  then return LastC
18                  else  LastScore ← Score
19                        LastC ← c
```
---

Subgraph isomorphism is a NP-complete problem, the *Persona Mapping Algorithm* avoids this pitfall. It builds an intersection graph between two subgraphs greedily to check for similarity making the problem tractable. The time complexity of the algorithm is $O(n^4 + n \log n)$. However, real world social networks have node degrees which follow a power law distribution [2]. Therefore the edge mapping routine will take on average $O(d^2)$, where $d$ is the average node degree of the network. The use of $\Delta$ to limit the depth of the breadth first search results in a further decrease in running time to build the auxiliary graph. Empirically a depth of 2 or 3 is sufficient to achieve a correct mapping as larger depths will include the entire graph due to the small world phenomenon [8], which in most cases is unnecessary.

Percentage of personas mapped correctly is a good measure of performance on a given data set. It is argued by [10] that this metric is flawed because nodes which are impossible to map will bias the results. However, this does not matter while measuring relative performance on a given data set. Finding a mapping between one entity and its personas is as important as the other, using measures such as node degree or centrality to give more weight to the successful mapping of more important nodes is meaningless here. For the same reason doing well on personas with relatively less information associated with them is also not considered.

## 4 Experiments

### 4.1 Outline of Experiments

This section presents experiments which validate the *Entity-Persona Model* and the *Persona Mapping Algorithm* against the Enron data set. The experiments test the robustness of the algorithm and the semantic similarity measures it uses by artificially partitioning the data set. At least 50% of the entities in the data set had their personas correctly mapped in every case.

The Enron data set was initially released by the Federal Energy Regulatory Commission as part of its Western Energy Markets investigation. The particular version used by this work has been prepared by Shetty and Adibi [12]. This data set is an appropriate test bed for this research as it is only composed of text and the link structure is easy to infer. As it is a very mature data set, it required minimal preprocessing to make it usable. Its size also makes it appropriate for research while still exhibiting many of the properties of large social networks.

The Enron data set consists of 517,431 emails from the mail accounts of 151 Enron employees between January 1998 to December 2002 with most of the email volume occurring between January 2000 and December 2001. It contains emails that originate or are sent to addresses outside those of the employees in the data set, these have not been considered in these experiments. The *Entity-*
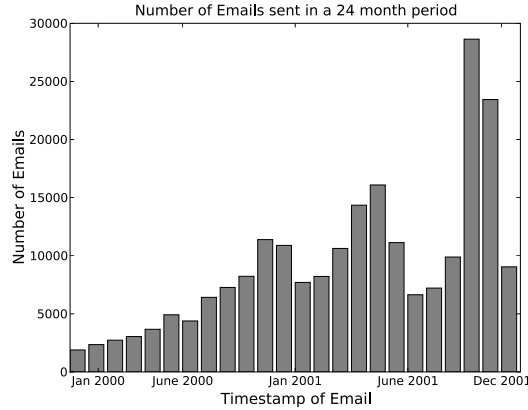


**Fig. 2.** Distribution of Enron email volume over a 2 year period.

*Persona Model* is populated by considering every email address as a node and a edge between two nodes is added if they have exchanged at least one email. The vertex labels of the social network are inferred from an unordered collection of terms in the text of all the emails sent or received by a node. The edge labels

103

are inferred from the common emails exchanged. Only the subject and body of the email is considered. To simulate anonymous interaction personal names and email addresses are filtered out. No stop words or stemming is used as this will not allow the stylometric features of the text to be captured. The goal of the persona mapping algorithm is to map personas unknown to it. A true test of its capabilities is its ability to map personas attempting to be anonymous. Although the web is full of this kind of data, it will be impossible to verify if the algorithm has made a correct mapping as this would require entities to volunteer this information. The next best approach is to synthetically create data sets to run the algorithm on, allowing the mappings to be easily verified. This research uses *temporal partitioning* on the Enron data set.

The data set is partitioned into eight sets based on when an email is sent as shown in table 1. The effect of the change in topics of the emails and change in graph structure can be explored. From table 2 it can be seen that each of the partitions has a different graph structure. At the local view of a node the graph changes significantly with some nodes communicating with completely new neighbors. The partition with the largest email volume is chosen to derive the source graph, this is the graph which the mapping between an entity and a persona is known. The seven other partitions are used to derive the target graphs whose personas need to be mapped.

|  |  | Partition by Email Date | |
| --- | --- | --- | --- |
| Graph Type | Graph Name | Start Date | End Date |
| Source Graph | $G_S$ | $1^{st}$ October 2001 | $31^{st}$ December 2001 |
|  | $G_{T1}$ | $1^{st}$ July 2001 | $31^{st}$ August 2001 |
|  | $G_{T2}$ | $1^{st}$ April 2001 | $30^{th}$ June 2001 |
|  | $G_{T3}$ | $1^{st}$ January 2001 | $31^{st}$ March 2001 |
| Target Graph | $G_{T4}$ | $1^{st}$ October 2000 | $31^{st}$ December 2000 |
|  | $G_{T5}$ | $1^{st}$ July 2000 | $31^{st}$ August 2000 |
|  | $G_{T6}$ | $1^{st}$ April 2000 | $30^{th}$ June 2000 |
|  | $G_{T7}$ | $1^{st}$ January 2000 | $31^{st}$ March 2000 |

**Table 1.** Partitioning the Enron Data Set.

### 4.2 Comparison of Semantic Similarity Techniques

The semantic similarity techniques presented in section 3.3 which are used by the *Persona Mapping Algorithm* are compared here. The experiment measures the percentage of personas successfully mapped between the source and seven target data sets. The success of these techniques depends on term usage by the entities in the the data sets. An analysis of the dataset showed that the number of unique words used by an entity is weakly correlated with the volume of communication.

| Statistic | $G_S$ | $G_{T1}$ | $G_{T2}$ | $G_{T3}$ | $G_{T4}$ | $G_{T5}$ | $G_{T6}$ | $G_{T7}$ |
|---|---|---|---|---|---|---|---|---|
| Nodes | 139 | 137 | 140 | 109 | 107 | 93 | 79 | 58 |
| Node Degree (Max) | 63 | 54 | 63 | 24 | 22 | 18 | 11 | 14 |
| Node Degree (Avg) | 11 | 8 | 7 | 6 | 5 | 4 | 3 | 2 |
| Edges | 825 | 616 | 531 | 345 | 328 | 234 | 127 | 88 |
| Email Volume (Avg) | 531 | 466 | 418 | 457 | 448 | 432 | 306 | 338 |
| Email Volume (Med) | 257 | 192 | 116 | 66 | 58 | 24 | 2 | 0 |

**Table 2.** Graph Statistics

This is the basis of the success of using textual information to find mappings as an entity does not have to have communicated extensively for there to be enough information to uniquely identify it.

The algorithm was run with $\sigma$ set as 0.6, $\Delta$ as 1 and $\epsilon$ as 0.002. These settings give more weight to the score of the mapping of the candidate vertex to the source and less to the scores of their common neighbors. This can interpreted as the algorithm has less confidence in that entities use their personas to communicate with the same people and more confidence in that exchanging similar information across all their personas. A smaller $\Delta$ implies the algorithm only looks at the immediate neighbors of the personas for structural congruence.
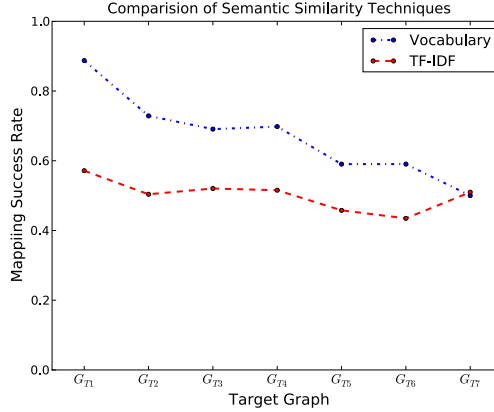


**Fig. 3.** Comparison of Semantic Similarity Techniques.

The Jaccard index to compare persona vocabulary was the more successful semantic similarity measure. It was also the least consistent in producing successful mappings as the time between which an entity used its persona increased. This implies it is the least robust to changes in specific topics of communication.

105

| Semantic Similarity Technique | Average Mapping Time |
| --- | --- |
| Jaccard Index | $2.18s$ |
| TF-IDF | $25.25s$ |

**Table 3.** Comparison of Time Required to Map

It is also the least computationally intensive technique to transform raw text into a representation which can be compared. Only one pass was required to populate the vocabulary vectors which label the graph.

TF-IDF performs more consistently although with a lower correct mapping rate. As term frequency is taken into account, it is less sensitive to rarely used terms. It correctly mapped the same subset of entities through all the data sets. It requires at least two passes over the email text to form the TF-IDF vector and is very slow at finding a mapping because the length of its vector is the size of the vocabulary of all the personas.

## 5 Conclusions

This research formally defines the problem of identifying the personas of entities on the web and proposes a solution to identify the two personas of an entity which exist in two separate social networks. This solution can be easily extended to the more general problem of multiple social networks and multiple personas. This research aims to answer the question: "Who am I interacting with?" or in other words, "Can I trust the identity presented as a true representation of the entitys identity?" This research takes a first step by offering a measurement of similarity among digital personas to resolve the true identities. Consequently, increased assurance of digital identities will allow for increased trustworthiness for on-line interactions with unknown entities (e.g. email, social networks, e-commerce, etc.).

The *Entity-Persona Model* is a formal model of the social graph and the *Persona Mapping Algorithm* operates on this model to produce mappings between entities and their personas. The *Entity-Persona Model* was populated using the Enron data set and the performance of the proposed algorithm was studied. Although the algorithm was successful in correctly mapping the entities in the data set to their personas, it was not consistent in its performance. The vocabulary, TF-IDF, and topic distribution approaches for semantic similarity compared in this research differ in their robustness to change in topics of information exchanged and the ability to capture stylometric information. These approaches can be improved or combined to perform consistently across all personas.

The question of what information generated by an entity is necessary and sufficient to uniquely identify it is the hardest to answer. It is apparent that the relative information entropy between the digital signatures of personas to be compared must differ by at least one bit in order to discriminate between them, which gives a definite lower bound for necessary information and can occur in a

fully connected social network. This lower bound however is not very useful in a real world situation as every virtual persona will have a unique signature. An empirical bound on the sufficiency of information can be found if all personas in the data set have been mapped correctly. However it will be impossible to know if all the personas have been mapped correctly as the purpose of the algorithm is to find unknown mappings and there is no ground truth to reference. User feedback can be used to learn how much information is sufficient by asking the user if he is satisfied with the results of the mapping, and then adjusting the amount of information considered based on this feedback. The danger of additional information adding more noise than signal must be considered while using such techniques.

## References

1. Ro Acquisti and Ralph Gross. Imagined communities: Awareness, information sharing, and privacy on the facebook. In *In 6th Workshop on Privacy Enhancing Technologies*, pages 36–58, 2006.

2. Lada A. Adamic, Rajan M. Lukose, Amit R. Puniyani, and Bernardo A. Huberman. Search in power-law networks. *Phys. Rev. E*, 64(4):046135, Sep 2001.

3. Lars Backstrom, Cynthia Dwork, and Jon Kleinberg. Wherefore art thou r3579x?: anonymized social networks, hidden patterns, and structural steganography. In *WWW '07: Proceedings of the 16th international conference on World Wide Web*, pages 181–190, New York, NY, USA, 2007. ACM.

4. L.P. Cordella, P. Foggia, C. Sansone, and M. Vento. A (sub)graph isomorphism algorithm for matching large graphs. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 26(10):1367–1372, Oct. 2004.

5. Judith Donath. *Identity and Deception in the Virtual Community*. Routledge, London, 1999.

6. James Grimmelmann. Facebook and the social dynamics of privacy. Draft article, August 2008.

7. Yingzi Jin, Yutaka Matsuo, and Mitsuru Ishizuka. Extracting social networks among various entities on the web. In *ESWC '07: Proceedings of the 4th European conference on The Semantic Web*, pages 251–266, Berlin, Heidelberg, 2007. Springer-Verlag.

8. Jon Kleinberg. The small-world phenomenon: An algorithmic perspective. In *in Proceedings of the 32nd ACM Symposium on Theory of Computing*, pages 163–170, 2000.

9. Ravi Kumar, Jasmine Novak, Prabhakar Raghavan, and Andrew Tomkins. Structure and evolution of blogspace. *Commun. ACM*, 47(12):35–39, 2004.

10. Arvind Narayanan and Vitaly Shmatikov. De-anonymizing social networks. Mar 2009.

11. Jasmine Novak, Prabhakar Raghavan, and Andrew Tomkins. Anti-aliasing on the web. In *WWW '04: Proceedings of the 13th international conference on World Wide Web*, pages 30–39, New York, NY, USA, 2004. ACM.

12. Jitesh Shetty and Jafar Adibi. The enron email dataset database schema and brief statistical report. Technical report, Information Sciences Institute, 2004.

13. Jessica Staddon, Philippe Golle, and Bryce Zimny. Web-based inference detection. In *SS'07: Proceedings of 16th USENIX Security Symposium on USENIX Security Symposium*, pages 1–16, Berkeley, CA, USA, 2007. USENIX Association.

14. Sherry Turkle. *Life on the Screen: Identity in the Age of the Internet.* Simon & Schuster, September 1997.
15. S. Umeyama. An eigendecomposition approach to weighted graph matching problems. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 10(5):695–703, Sep 1988.
16. Vassilios S. Verykios, Elisa Bertino, Igor Nai Fovino, Loredana Parasiliti Provenza, Yucel Saygin, and Yannis Theodoridis. State-of-the-art in privacy preserving data mining. *SIGMOD Record*, 33:2004, 2004.