

# A Model for Identity in the Cyber and Natural Universes

Duncan Hodges, Sadie Creese and Michael Goldsmith  
Cyber Security Centre,  
Department of Computer Science,  
University of Oxford, Oxford, UK

Email: duncan.hodges@cs.ox.ac.uk, sadie.creese@cs.ox.ac.uk, michael.goldsmith@cs.ox.ac.uk

**Abstract**—In this paper we describe a model for evaluating and investigating holistic identities across both the cyber- and natural-worlds. The model is taxonomy-agnostic and helps guide analysts through the creation of a rich superidentity whilst also being flexible enough for users to deploy their intuition and experience. Superidentities are created by iterative stages of enrichment and relaxation, which allow the superidentity to grow and include more elements of identity.

The confidence in individual elements can be propagated so as to provide an overall level of confidence in the superidentity. The quality of the final, rich superidentity can be metricized, identifying good qualities of the superidentity such as reinforcement of elements and consensus within element types.

This work is a key enabling component of a large interdisciplinary project called Superidentity, which is seeking to develop the model through understanding the elements of identity which might be available and how their reliability might vary according to context. In addition to this aim the project looks to develop visual analytic support to help validate the modelling approach. In particular demonstrating the model's use by analysts seeking to make identity attributions.

## I. INTRODUCTION

Modern society sprawls across cyberspace, engaging with it to obtain entertainment, education and emotional support [1] as well as a mechanism for purchasing products and services. As society's appetite for digital technology shows no sign of abating, an individual's exposure to cyberspace will continue to grow. Mobile technologies continue to provide new access opportunities allowing a rich 24-hour portable access experience. This is eloquently demonstrated by the activation rate of Android devices which is estimated to exceed 700,000 devices per day [2] which is almost twice the estimated growth-rate of the human species.

Cyberspace, and digital communication in general, is changing how we interact, consume media and define ourselves. The explosion of web 2.0 has created an environment such that, rather than passively consuming media, users actively provide new content, remix or 'mash-up' others' content as well as interacting with each other to comment and discuss content. It is perhaps not surprising that communities exploit this powerful tool to recruit new members [3]–[5], and there is significant evidence that a variety of extremist groups actively exploit cyberspace for fundraising [6] and as an operational tool (e.g. for attack planning and preparation).

Cyberspace is also a platform for more general criminal activity [7]. It is perceived to possess several key characteristics which make it an attractive platform for such activity, including: the perceived level of anonymity; its massive reach (meaning a large number of people can be targeted); and a low-level of operational risk (whether from law-enforcement or other criminals); whilst the victim is typically depersonalised potentially reducing the sense of guilt associated with crime [8]. In addition to these perceptions, the technical level of entry into cyber-crime is relatively low.

Elements of an individual's cyber-identity may be spread widely across cyberspace – distributed through Online Social Networks (OSN), blogs, web-fora or other content (e.g. photo sharing websites), online retailers and recommendation systems, etc. There are also channels where a user may be unaware that they are sharing information (e.g. through metadata associated with content, such as a geo-tagged image) or their footprint can be aggregated in order to gain author identification and other intelligence (e.g. [9]). It is this unconscious sharing which is difficult to mask and so may be one of the keys to exploiting the full identity exposure.

Other elements of identity exist or are derivable in the natural world, such as biometric measurements (e.g. fingerprints, vein patterns [10], ear patterns [11]). These are all indicative of an identity, and may cross the cyber / natural world divide through connection to cyber identities (such as associations with real-names which may also appear in email addresses or usernames).

The working hypothesis of this project is that there is a model of identity within which we can connect elements from both cyber and physical space, and that through the construction of such a model we can significantly enhance our capability to make identity attributions and facilitating a more holistic view that can be measured and assigned a level of confidence within a given context. The union of the elements of cyber-identity and the elements of natural-identity we call a **superidentity**, and by fully enriching a superidentity across a wide aperture we expect that issues such as deception or inconsistencies will appear as conflicts within the superidentity.

Identity, in all of its forms, in cyberspace is of interest to a number of different communities. The first of these communities comprises the users themselves; one of the core

concepts of social networks is the ability to interact with other users. In order to interact and build online relationships it is important to be able to identify users who are either friends or acquaintances in the real-world or other users that share similar interests. In order to achieve this it is important that services have an ability to identify and understand some concepts of identity relating to that particular service.

The next community which is interested in identity, and indeed is one of the key drivers for the recent explosion of identity in cyber-space, is primarily concerned with monetising ‘identity’. The largest commercial interest is undoubtedly advertising; the ability to profile a given user’s identity in terms of gender, recent searches, age, etc., drastically increases the success of targeted advertising and hence improves the financial return on the screen real-estate.

In addition to this community, there is an equally important educational-awareness facet: as the public exposes more information within cyberspace, the available vectors for social-engineering attacks and identity theft increase. So it is critically important that, as society increases the amount of its identities they share, the ramifications of this sharing should be understood.

One of the most important uses of identity within cyberspace is law-enforcement and intelligence provision. This largely takes the form of two types of activity. The provision of forensic evidence is important, as miscreants use digital communication for a host of illegal activities, it is important that law-enforcement has the skills and tools necessary to be able to attribute behaviours to identity. The act of attributing these behaviours to identity is generally associated with the provision of evidence for legal proceedings. The next type of activity is intelligence-gathering, where the investigator is trying to facilitate disruption of an undesired behaviour. In this case, the aim of the investigator is either to identify an individual from an MO-trigger or to attribute behaviours to given a target identity.

In this paper we present a mathematical model for superidentity, which supports exploration of identity and its multiple facets, the growth of a superidentity, the propagation of confidence in elements across a superidentity, and the measurement of superidentity quality. The construction of the model is a key enabling component of a large inter-disciplinary consortium project, itself called Superidentity [12], [13], which is seeking to develop the model through understanding the elements of identity that might be available and how their reliability might vary according to context. In addition visual analytic support will be used to validate the model’s utility for use by analysts seeking to make identity attributions, particularly those involving entities in cyberspace.

The model defined here explicitly allows superidentities to be created and explored, and superidentities themselves can be merged and split as intelligence is added and it becomes apparent that a single superidentity relates to a number of individuals, or what was thought to be a number of superidentities actually pertain to just one. The authors understand that a large part of intelligence-gathering exploits the intuition and

experience provided by the human analysts, so the framework is designed to be flexible enough to enable this human factor to be projected into the superidentity.

It is worth noting that this work is concerned with a subtly different problem than often seen in fields such as Digital Identity Management (e.g. [14]) where the user is typically compliant and provides (or is influenced to provide) high-quality, non-ambiguous elements of identity, and the ultimate aim is to completely automate the associations. In contrast we are concerned with environments where each of the elements is likely to be noisy, uncertain and commonly deliberately deceptive and, whilst automation is certainly desirable to support many simple tasks, the intuition and reasoning power of experienced analysts will remain essential.

The remainder of this paper is organised as follows: In II SuperIdentities and elements of Identity section.2 we describe how a superidentity is constructed from identity elements. In III The propagation of confidence within the elements of identity section.3 we present a strategy for propagating confidence through the superidentity. In IV Grouping of elements of SuperIdentities section.4 we describe how elements can be grouped to support the connection of identity threads to a single superidentity. In V Quality of SuperIdentity section.5 we present measures of quality in a superidentity. In VI New elements of identity section.6 we describe how new identity elements can be added to a superidentity. Finally in VII Conclusion section.7 we present conclusions and in VIII Further work section.8 future work.

## II. SUPERIDENTITIES AND ELEMENTS OF IDENTITY

An element of identity is defined as a single piece of information which can be either telegraphic, or indicative of an identity. Example elements of identity are biographic natural-world elements such as real-name, date-of-birth, home-address, occupation, etc.. or cyber identifiers such as usernames on a given website / service, email-addresses, etc. Other elements are more complex for example content on a blog or forum, location information, CCTV imagery, etc.

In this environment of noisy, deceptive (and indeed often conflicting) elements one must be cautious not to rush to formal ontologies (such as presented in [15]). By describing these ontologies in very formal description languages, such as OWL, it is possible to provide a very rigorous framework (which is desirable). However, this is often at the expense of flexibility; flexibility is particularly important if the system is to fully exploit the human analyst.

Mathematically, an element of identity of type  $a$  is denoted by  $e_a$ :  $a$ . A characteristic is a multiset of elements of identity of the same type, denoted by  $C$ ; whilst a superidentity,  $S$ , is a set of these characteristics (and hence contains the elements). Note that the SuperIdentity is **not** a multiset since it only contains a single set of each type of characteristic. For example, a superidentity consisting of two elements of type  $a$ , a single element of type  $b$  and an element of type  $c$  can be

written as :

$$C_a : a = \{e_a : a, e'_a : a\} \quad (1)$$

$$C_b : b = \{e_b : b\} \quad (2)$$

$$C_c : c = \{e_c : c\} \quad (3)$$

$$S = \{C_a : a, C_b : b, C_c : c\} \quad (4)$$

Constructing a superidentity must start with at least one *seed* element of identity, e.g.  $e : a$  (this could be a biographical identifier such as a real-name, or a cyber identifier such as an email address, Twitter username or Facebook ID), this forms the initial superidentity. The superidentity is then enriched by taking each of the elements of identity and deriving new elements of identity; for example an email address may be transformed to produce usernames on social network sites, company names, etc. in addition to the original email address.

Essentially there is a collection of transform functions that transform an element of a particular type to another type e.g. for a source element of type  $a$  there is a transform  $t_{a \rightarrow b}$  that creates an element of type  $b$  from an element of type  $a$  such that:

$$t_{a \rightarrow b} : e_a : a \mapsto e_b : b \quad (5)$$

or to take a tangible example:

$$\begin{aligned} & t_{EmailAddress \rightarrow FlickrUsername} : \\ & \text{alice@yahoo.com} : EmailAddress \mapsto \\ & \text{alicephotographer} : FlickrUsername \end{aligned} \quad (6)$$

Where  $t_{EmailAddress \rightarrow FlickrUsername}$  is a method in the Flickr API [16]. The function for the example in (5SuperIdentities and elements of Identityequation.2.5) shows a trivial one-to-one transform. There are also one-to-many transforms (such as mapping a real-name to a Facebook account<sup>1</sup>), many-to-one transforms (such as querying birth records which require the real-name, place-of-birth and potentially the date-of-birth) and many-to-many transforms (such as real-name and employer to email-address).

Many-to-many or one-to-many transforms where there is a definable set of outputs can be modelled as a number of many-to-one or one-to-one transforms. For example the transform that uses the real-name and employer to create an email-address can be modelled as a number of different transforms: one which produces first-name.surname@employer, another produces first-initial.surname@employer, another produces first-name@employer etc. This has the advantage that each of these transforms can have a different confidence (for more discussions on the confidences see IIIThe propagation of confidence within the elements of identitysection.3).

For example the many-to-one transform that creates an element of type  $d$  from elements of type  $b$  and  $c$  can be described:

$$t_{(b \times c) \rightarrow d} : (e_b : b \times e_c : c) \mapsto e_d : d \quad (7)$$

<sup>1</sup>and indeed the set returned may not include the account associated with the target

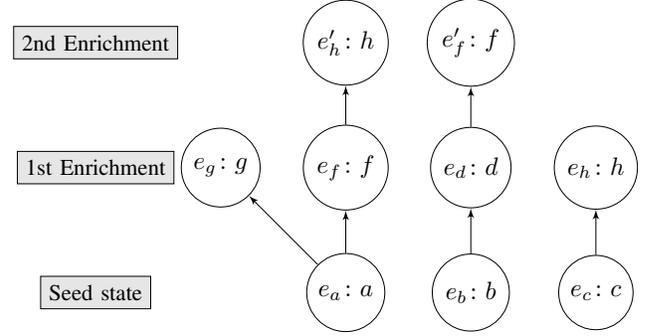


Fig. 1. Example superidentity from three seeds after two levels of enrichment

or to take a example:

$$\begin{aligned} & t_{(RealName \times Employer) \rightarrow EmailAddress} : \\ & (\text{BobCharles} : RealName \times \\ & \text{DodgyElectronics} : Employer) \mapsto \\ & \text{bob.charles@dodgyelectronics.com} : EmailAddress \end{aligned} \quad (8)$$

If the set of outputs is not definable at the outset (e.g. mapping a real-name to a Facebook username) then a single transform is implemented which produces a set of output elements which have equal confidence. For example if there exists a transform which creates a set of elements of type  $g$  from an element of type  $f$ :

$$t_{f \rightarrow g} : e_f : f \mapsto \{e_g : g, e'_g : g, e''_g : g\} \quad (9)$$

The concept of the *reachability matrix* has been developed in [17] and is an area of ongoing research. The framework presented here will sit beside this work and create a powerful ensemble.

The process of enriching the superidentity continues by iteratively transforming the elements which form the superidentity, this creates a directed graph which outlines the provenance of the elements of identity, an example graph from three seeds is shown in Fig.1Example superidentity from three seeds after two levels of enrichmentfigure.1. It is clear that the number of elements per enrichment shown in Fig.1Example superidentity from three seeds after two levels of enrichmentfigure.1 is relatively small, this is for the sake of clarity – real enrichments are likely to develop many more elements.

### III. THE PROPAGATION OF CONFIDENCE WITHIN THE ELEMENTS OF IDENTITY

Every element of identity has a confidence,  $R$ , associated with it such that  $0 \leq R \leq 1$ , where a confidence of 1 represents the fact this element is ‘correct’ and is an element associated with the superidentity; 0 represents an element which is not possibly part of the superidentity or is known to be false. For example the element  $e_b : b$  derived in (5SuperIdentities and elements of Identityequation.2.5) has a confidence  $R(e_b : b)$  or  $R(\text{alice@yahoo.com} : EmailAddress)$ . The value  $R$  can either be a single discrete value or a probability density

function (PDF) describing the likely confidence of a range within the element (for example a date-of-birth element of identity could itself be a PDF rather than a single discrete value – it is worth noting that the number of elements which are of a continuous nature is likely to be relatively small compared to those elements which are discrete in nature).

The confidence of the derived element from a functional transform can be calculated by:

$$R(e_b : b) = R(e_a : a) \cdot R(t_{a \rightarrow b}) \quad (10)$$

The confidence of a particular element is clearly related to the confidence of the element from which it is derived and the confidence in the enrichment process. In the example enrichment from email address to Flickr username it is clear that pragmatically the enrichment cannot make an error<sup>2</sup>. If this simple enrichment cannot be in error then the confidence in the Flickr username should be that of the source element in this case the confidence in the transform is high  $R(t_{EmailAddress \rightarrow FlickrUsername}) = 1$ .

In the more complex case where there is a compound input, the confidence among the input elements is combined with the confidence of the transform, such that considering the example in (7SuperIdentities and elements of Identityequation.2.7) the new confidence is calculated:

$$R(e_d : d) = R(e_b : b) \cdot R(e_c : c) \cdot R(t_{(b \times c) \rightarrow d}) \quad (11)$$

It is worth noting that in some cases the enrichment confidence,  $R(t_{a \rightarrow b})$  will not be a constant but will be a function of the input element  $e_a : a$ . This is an important observation, for example a CCTV image of a suspect may have a very high confidence (i.e. the analyst is very sure this is the individual concerned) and there will be a transform which attempts to map this image to a gender. The confidence of this transform is not solely a function of the algorithm but also it is dependent on the quality of the image.

As can be seen the confidence in a particular element is very strongly tied to the elements provenance, the directed graph formed by the elements (vertices) and the transforms (edges) can be used to intuitively retain this information.

#### IV. GROUPING OF ELEMENTS OF SUPERIDENTITIES

The definition of the superidentity as a set of characteristics can be emphasised by redrawing the example superidentity shown in Fig.1Example superidentity from three seeds after two levels of enrichmentfigure.1 is redrawn in Fig.2Example superidentity rearranged to highlight elements of identical typefigure.2 to highlight elements that are the same type.

Within a superidentity it is common that some elements of identity of the same type may be produced from different sources, for example an employer may be acquired from OSNs, email address, location traces, etc. and this forms the characteristic *Employer*. This multiset contains the various estimates of this particular type from the superidentity, each

<sup>2</sup>The concepts of deliberate deception by the target e.g. the use of a pseudonym is covered later, the only way the returned value can not be ‘correct’ is direct malicious interference with the API call

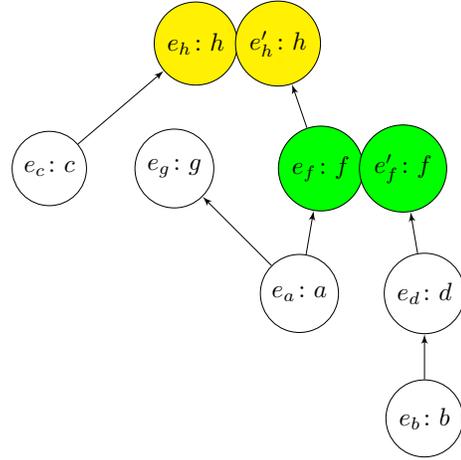


Fig. 2. Example superidentity rearranged to highlight elements of identical type.

of the elements within the multiset has a confidence associated the element’s provenance.

For each characteristic it should be possible to define a similarity function,  $L$ , that measures the similarity between two instances of the same type. Taking the example shown in Fig.2Example superidentity rearranged to highlight elements of identical typefigure.2 there should be a measure for element type  $f$  such that if  $l = L_f(e_f : f, e_f\' : f)$ ,  $0 \leq l \leq 1$ . This function will vary depending on the type of element: if the elements are real-names the similarity of two elements can be measured by the edit-distance, if it is two profile pictures then facial recognition techniques might be used [18], [19]; if it is two location traces (e.g. one from mobile telephony, one from Google Latitude or other geo-tagged metadata) then notions such as *unlikely-* or *impossible-travel* can be used.

In general a mechanised approach will provide only an estimate of the similarity as a decision-enhancing tool, and ensuring a human provides insight will allow a much more complete assessment of the similarity. Indeed the analyst may override the mechanised function  $L$  to assign a similarity  $l'$ .

It is now possible to consider a group of elements of the same type: each of the elements has an estimate of its similarity to all the other elements of the same type. This grouping provides a mechanism through which the particular characteristic can be analysed; elements, whether they are similar or not, retain their provenance as they are grouped and hence retain their confidence measures.

If the similarity between elements is high but not exactly one (e.g. there are small differences between elements) but the analyst has overruled the machine-based similarity function (such as where one element forms the superset of another e.g. work-address is Oxford and work-address is Keble Road, Oxford) then the analyst has, effectively, imparted some external knowledge that these two facts are consistent and set the similarity function to  $l' = 1$ .

The combination of the confidence of individual elements

within the characteristic and the similarity amongst all the elements provides a good insight into both the consistency and reinforcement of that particular characteristic. For example, if a group has a number of elements from different sources all with a low confidence but they are all similar then it may mean that the confidence in this set of elements is actually higher.

In essence the model retains the elements when they are grouped allowing individual elements to retain their provenance (and hence ‘confidence flow’): as characteristics are grouped, then confidence in the elements may increase and this increase is fed back to the downstream<sup>3</sup> nodes. These dynamic changes in the confidence propagation provide a system such that as it is enriched and elements are reinforced, there is a rebalancing of confidence gravitating towards a stable system where characteristics have a high confidence and low conflict.

To this point we have only considered characteristics where the elements are similar within a characteristic, of course in a number of situations it is possible to have elements which are in conflict. Understanding and using this conflict is important and will depend on the identity-attribution exercise that is being pursued. If, for example, the exercise is investigating a single individual and the values of  $l$  or  $l'$  are small, the analyst cannot confidently say the two elements are consistent. In this case there is conflict in this element type within the characteristic, this can occur in a number of situations.

The first circumstance in which this occurs is when a target is being deceptive – this may be conscious, deliberate deception such as a target working under multiple personae or more subtle ‘deception’ deployed in order to protect a user’s perception of privacy and safety [20], [21]. Interestingly, it maybe impossible to detect this deception without taking a holistic view of the superidentity, particularly with highly OPSEC-aware ‘hard’ targets seen in some intelligence or law-enforcement environments.

The second such circumstance is similar to the first and occurs when there genuinely are multiple elements of the same type (e.g. a single user may have multiple email accounts). This is subtly different to the first case, in this case the elements are all valid whilst in the earlier case there is only one valid element. These two cases will be difficult to distinguish, and we can incorporate these ‘deceptive’ personae within the superidentity, for example in order to maintain information pertaining to persistent and consistent pseudonyms.

The third circumstance in which the elements of identity may be in conflict is where there are elements which should not be part of the superidentity. This is most likely to arise during an investigation where a seed element is not globally unique. If a seed element is globally unique there are likely to only be functional transforms out to new elements; e.g., an email address can only be linked to one Facebook username. Transforms from non-unique seeds are more likely to be one-to-many mappings, hence a real-name may map to multiple OSN usernames or have multiple entries in phone-books.

<sup>3</sup>We refer to a node  $i$  being downstream of node  $j$  if there is a path leading from  $j$  to  $i$  and no path leading from  $i$  to  $j$ .

The final condition is where one enrichment is incorrect, for example an element may not be of a sufficiently high intrinsic quality<sup>4</sup> to generate a correct enrichment. It is anticipated that the confidence of the transform will be low and hence the derived element will be of a low confidence.

Before an element is grouped it possess a confidence purely from it’s provenance,  $R(e_a : a)$ . However, when the element is grouped within a characteristic then the confidence of the element should take account of the confidence and proximity of the other elements within the characteristic. In this parlance, proximity refers to the similarity between the elements – i.e. it’s proximity in the topology of the characteristic.

We can define a measure which is a measure of both the confidence derived from the provenance of the element **and** the proximity to the other elements in the characteristic such that:

$$R_a^\dagger(e_a : a) = \phi \cdot R(e_a : a) + (1 - \phi) \cdot \left[ \sum_{e'_a : a \in C_a} R(e'_a : a) \cdot L(e_a : a, e'_a : a) \right] \quad (12)$$

The value of  $\phi$  defines how much this characteristic is affected by the similarity of other elements of the same type and varies between element types. The measure takes the confidence from the provenance and adds a contribution from the other element confidences in the characteristic (modulated by the similarity between the two elements).

## V. QUALITY OF SUPERIDENTITY

As can be seen the superidentity is built-up following successive enrichment and grouping allowing the superidentity to grow and provide a rich complete picture of the individual. In essence the directed graph of the elements of identity can be imagined as a mechanical system with the elements of identity as nodes and the enrichments representing the edges. As the superidentity is enriched, the system is energised before relaxing as elements are grouped. The model allows the analyst to input their own confidence values and decouple the feedback of confidence to up-stream elements. For example if two elements of type ‘gender’ are grouped, one of which is very high confidence and one of which is very low confidence it may be undesirable to feedback the increase in confidence.

There is value in being able to metricize the quality of the superidentity, in addition to providing guidance as to what new seed elements would provide the most valuable opportunities to discover new elements. The second point can be achieved through analysis of the confidence / provenance within the directed graph in addition to the reachability matrix, desirable elements can be identified in order to prioritise new seed elements (so as to provide guidance for future humint operations).

The first point, that of measuring the quality of a superidentity, is more abstract and putatively difficult to metricize.

<sup>4</sup>It should be noted that it is possible for an element to have a high confidence but a low intrinsic quality - such as low resolution images or voice-cuts

Essentially a good superidentity should be rich (i.e. have a large number of characteristics), demonstrate consensus and reinforcement.

The reinforcement can be best visualised from the directed graph of the elements, a directed graph which is a tree (in which any two nodes are connected by exactly one vertices), or more likely a forest (the disjoint union of trees) represents a superidentity with no reinforcement and cannot be described as notionally ‘good’. Whilst consensus is the degree to which the various elements within a characteristic are similar, a better way of considering consensus is the lack of conflict. Some characteristics may only have one valid element (for example an individual can only have one date-of-birth) whilst some characteristics may have multiple valid elements which can co-exist within a superidentity without conflict, for example, location at various times. However, these elements can come into conflict if travel between two points is unlikely.

As can be seen consensus and reinforcement are related whilst being subtly different, for example with no reinforcement a superidentity will have perfect consensus. Equally a small superidentity may have good consensus and reinforcement but will not be a rich representation of the superidentity. A superidentity can only be considered notionally ‘good’ if it possess all three features: consensus (or lack of conflict), reinforcement and richness (a large number of characteristics).

In order to attempt to metricise the quality of a particular characteristic we can use a definition of the form (13Quality of Superidentityequation.5.13) to express the effective diversity of a characteristic. For each element, we calculate a negative exponential of it’s ‘distance’ (defined in terms of the similarity,  $L$ ) from each other element, but modulate this by the confidence of the other point, so that low-confidence elements have less effect.

$$D_a = \frac{\sum_{x,y \in C_a: a} \left[ R_a(x: a) R_a(y: a) \exp\left(1 - \frac{1}{L(x: a, y: a)}\right) \right]}{\sum_{x,y \in C_a: a} [R_a(x: a) R_a(y: a)]} \quad (13)$$

A wholly self-corroborating characteristic  $C_a$  will have  $D_a = 1$ , whilst a  $C_a$  with equal confidence in entirely unlike elements will have  $D_a = 1/|C_a|$ .

This measure of the quality of a particular characteristic can be aggregated to produce an overall quality of the superidentity, organisations may weight the various characteristics to produce a score based on their requirements, or the capabilities to be deployed against particular targets. Heuristically, if two elements agree with the gender of a superidentity it *means* less than if two sources agree with the home-address. The *value* of the consensus will vary across the characteristics and will be affected by a number of factors, such as the uniqueness of the element (for example home-address is more of a unique element than gender so hence putatively a higher value) and the usefulness of the element (an element of identity which allowed targeting of an individual will be higher value than

one that does not<sup>5</sup>).

This implies that all characteristics of identity are not equally important, and the relative importance of individual elements will vary across organisational requirements and indeed may vary across identity-attribution exercises within one organisation. It is important to attempt to characterise this within the metrication of the superidentity, a superidentity with only low ‘value’ elements of identity should itself achieve a lower score than that of a superidentity which has higher ‘value’ elements.

In conclusion, (13Quality of Superidentityequation.5.13) is a measure of the quality of a superidentity which effectively metricises the correlation between elements weighted by their confidence, whilst the individual confidence measures associated with each element provides a higher resolution fidelity which can be used identify weak and strong elements within the superidentity.

## VI. NEW ELEMENTS OF IDENTITY

In order to introduce a new element of identity into the framework there is a simple set of requirements. Initially the enrichment processes from that element should be created (in essence defining how to derive elements from this new element type),  $T_{\text{new}}$ . Enrichment transforms should also be added to any existing types which can be used to derive this new element type, e.g.  $t_{a \rightarrow \text{new}}$ . This should include the functions that are able to calculate the confidence in the new elements, e.g.  $C(t_{a \rightarrow \text{new}})$ .

These functions effectively define the relationship between the new element type and the old element types. The reachability matrix,  $T$ , will now include new information so the transitive closure should be recalculated as the ‘range’ of seed elements will now be changed.

Interestingly, the framework presented in this paper could provide a way of automatically handling ‘cold-cases’: the superidentity can be archived, and as new capability is discovered or invented the reachability matrix will change and hence there is a possibility of new intelligence being automatically derived, and analysts notified. To take the example shown in Fig.1Example superidentity from three seeds after two levels of enrichmentfigure.1, if the derived superidentity is not sufficient to provide valued intelligence this superidentity may be archived; if in the future the capability associated with the enrichments  $t_{g \rightarrow q}$  and  $t_{q \rightarrow b}$  are produced, then it would be possible to automatically derive a new element  $e_q: q$  and reinforce the element  $e_b: b$ , which may significantly improve the superidentity and make it more valuable.

One of the hardest functions to define for a new element type is likely to be the similarity function,  $L$ , which is a measure of the degree to which two elements of same type are alike. This is required in order to assist the analyst with the merging stage of the superidentity enrichment.

These set of requirements should allow for the creation of a particular taxonomy of elements quickly and rigorously. The

<sup>5</sup>This may also be related to the uniqueness, i.e. unique elements will tend to be used for targeting.

taxonomy can be created to be dependent on the analyst's requirements from the framework, the available enrichments and other external factors (such as the overarching legal framework). The framework could also be used to collate information from other collection systems in a rigorous yet agile framework.

## VII. CONCLUSION

Society is at a pivotal stage of development, explosive growth in a mobile rich access to cyberspace has helped fuel and sustain a constant 'always-on' society where identities in the cyber world and the natural world are blurred, and these worlds can no longer be treated as separate.

In this paper we have presented a model in which these identities can be modelled and explored in order to make identity attributions. The holistic approach to identity exploration enabled by this framework will be key to identifying and exploiting deception and 'hard-targets' in cyberspace, and we believe it has particular value to the intelligence and law-enforcement communities, as well as to those concerned with understanding the exposure to risk that users of cyberspace face.

The model builds up a number of elements of identity, where each element represents a single atomic element of identity - these can be biographical information (such as real-name, date-of-birth, home-address, etc.), natural-world elements (such as location traces, biometric measures, CCTV imagery, etc.), cyber identities (such as usernames on a given site, email-addresses, etc.) or associated with access technologies (such as IMEI, IMSI, MAC addresses, ISP, etc.).

A set of enrichments each map one element of identity to another. This process of enrichment allows a Superidentity to grow from a set of seed nodes (representing the initial knowledge of a target). The complete set of enrichments forms a reachability matrix. Following the enrichment process the superidentity undergoes a grouping process (or relaxation) which attempts to group elements of the same type. Hence, as this process of enrichment and relaxation continues, the superidentity grows and expands to become a rich view of the holistic identity of the target concerned.

The model also provides a mechanism for propagating confidence along the directed graph defining the provenance of elements of identity. The model can be overridden by human interaction, and indeed the intuitive capability of experienced analysts is important in helping to build the superidentities. Ongoing work is looking at surveying the processes and intuitive approaches that experienced analysts in law-enforcement and intelligence agencies deploy while performing identity attribution exercises.

Finally, the model is designed to be taxonomy-agnostic, meaning that the model is flexible enough to be able to provide a mechanism for users to incorporate their own collection architectures and association frameworks for individual elements of identity, and the model manages the merging, distillation and confidence propagation.

## VIII. FURTHER WORK

On-going work within the Cyber Security Centre is investigating taxonomies which can eloquently describe the elements of identity both in cyber- and the natural world and provide tangible demonstrators of the modelling approach. These demonstrators provide a 'play-pen' in which these superidentities can be created and explored. Early demonstrators show the value of several novel visualisation approaches which allow analysts to experiment with the confidence of elements and explore how the confidence of the superidentity varies.

The demonstrator also allows analysts to perform sensitivity analysis on the holistic superidentity, identifying those elements which are pivotal to creating the superidentity and how these can be reinforced. The reachability matrix can also help to direct the analyst in recommending new elements which would significantly improve the superidentity (either in terms of the 'richness' or the confidence of the superidentity). This analysis can provide a strong direction for other information gathering exercises.

An interesting assumption to date is that the elements of identity are assumed to be temporally static. This may be true for most elements which are biographic in nature (for example, real-name, date-of-birth); but others may change over time, for example in the medium-term an individual's home-address may change and in the longer-term an individual's biometric elements may change (for example fingerprints may change over time due to injury) [22]. Other elements of identity, such as location, will inherently change over time and hence in the future a temporal element should be introduced into the model.

As the boundaries between cyber- and natural world identities are becoming less well-defined it is becoming increasingly important that, in order to successfully be deceptive, deception needs to cover both off- and on-line arenas (a mid-twenties western European individual without an OSN projection is suspiciously rare, as is a company without a website). A further challenge is that the cyber-world is instantaneous and persistent, inconsistencies are immediately broadcast and last forever. This research on the holistic modelling of identity is key to enabling research on both establishing successful deception and in the recognition of deception in others.

## ACKNOWLEDGMENT

The work is performed under EPSRC grant EP/J004995/1. The SuperIdentity project is investigating the interactions between offline and online identities, the cross-disciplinary consortium ranges from innovative new biometric measures through to management of on-line identities. The authors would also like to thank our colleagues in the Cyber Security Centre at the University of Oxford, particularly Michael Goldsmith, Jason Nurse, Thomas Gibson-Robinson and Elizabeth Phillips whose early work developing a transitivity model for relating identity elements has been instrumental in developing the superidentity model.

## REFERENCES

- [1] H. White, E. McConnell, E. Clipp, L. Bynum, C. Teague, L. Navas, S. Craven, and H. Halbrecht, "Surfing the net in later life: A review of the literature and pilot study of computer use and quality of life," *Journal of Applied Gerontology*, vol. 18, no. 3, pp. 358–378, 1999. [Online]. Available: <http://jag.sagepub.com/content/18/3/358.abstract>
- [2] E. Schonfeld. (2011, December) Android phones pass 700,000 activations per day, approaching 250 million total. [Online]. Available: <http://techcrunch.com/2011/12/22/android-700000/>
- [3] M. Sageman, *Leaderless Jihad: Terror networks in the 21st Century*, 1st ed. University of Pennsylvania Press, 2008.
- [4] P. B. Gerstenfeld, D. R. Grant, and C.-P. Chiang, "Hate online: A content analysis of extremist internet sites," *Analyses of Social Issues and Public Policy*, vol. 3, no. 1, pp. 29–44, 2003. [Online]. Available: <http://dx.doi.org/10.1111/j.1530-2415.2003.00013.x>
- [5] J. Qin, Y. Zhou, and H. Chen, "A multi-region empirical study on the internet presence of global extremist organizations," *Information Systems Frontiers*, vol. 13, pp. 75–88, 2011, 10.1007/s10796-010-9277-6. [Online]. Available: <http://dx.doi.org/10.1007/s10796-010-9277-6>
- [6] S. D. Keene, "Terrorism and the internet: a double-edged sword," *Journal of Money Laundering Control*, vol. 14, no. 4, pp. 359–370, 2011.
- [7] D. S. Wall, *Crime and the Internet: Cybercrimes and Cyberfears*. Routledge, November 2001.
- [8] E. A. Kallman and J. P. Grillo, *Ethical Decision Making and Information Technology*. McGraw-Hill Education, 1996.
- [9] P. Juola, "Authorship attribution," *Found. Trends Inf. Retr.*, vol. 1, no. 3, pp. 233–334, Dec. 2006. [Online]. Available: <http://dx.doi.org/10.1561/1500000005>
- [10] J. Cross and C. Smith, "Thermographic imaging of the subcutaneous vascular network of the back of the hand for biometric identification," in *Security Technology, 1995. Proceedings. Institute of Electrical and Electronics Engineers 29th Annual 1995 International Carnahan Conference on*, oct 1995, pp. 20–35.
- [11] B. Moreno, A. Sanchez, and J. Velez, "On the use of outer ear images for personal identification in security applications," in *Security Technology, 1999. Proceedings. IEEE 33rd Annual 1999 International Carnahan Conference on*, 1999, pp. 469–476.
- [12] S. Black, S. Creese, R. Guest, B. Pike, S. Saxby, D. Stanton-Fraser, S. Stevenage, and M. Whitty, "Superidentity: Fusion of identity across real and cyber domains," in *ID360: The global forum on Identity*, April 2012.
- [13] SuperIdentity. (2012) Superidentity. [Online]. Available: <http://www.superidentity.org>
- [14] G. Ayed and S. Ghernaoui-Helie, "Digital identity management within networked information systems: From vertical silos view into horizontal user-supremacy processes management," in *Network-Based Information Systems (NBIS), 2011 14th International Conference on*, sept. 2011, pp. 98–103.
- [15] P. Rempel, B. Katt, and R. Breu, "Supporting role based provisioning with rules using OWL and F-logic," in *On the Move to Meaningful Internet Systems: OTM 2010*, ser. Lecture Notes in Computer Science, R. Meersman, T. Dillon, and P. Herrero, Eds. Springer Berlin / Heidelberg, 2010, vol. 6426, pp. 600–618. [Online]. Available: [http://dx.doi.org/10.1007/978-3-642-16934-2\\_45](http://dx.doi.org/10.1007/978-3-642-16934-2_45)
- [16] Flickr. (2012) Flickr API documentation. [Online]. Available: <http://www.flickr.com/services/api/flickr.people.findByEmail.html>
- [17] S. Creese, M. Goldsmith, J. Nurse, and E. Philips, "A data-reachability model for elucidating privacy and security risks related to the use of online social networks," in *International Workshop on Trust, Security and Privacy in e-Government, e-Systems and Social Networking*, June 2012.
- [18] M. Nechyba, L. Brandy, and H. Schneiderman, "Pittpatt face detection and tracking for the clear 2007 evaluation," in *Multimodal Technologies for Perception of Humans*, ser. Lecture Notes in Computer Science, R. Stiefelwagen, R. Bowers, and J. Fiscus, Eds. Springer Berlin / Heidelberg, 2008, vol. 4625, pp. 126–137. [Online]. Available: [http://dx.doi.org/10.1007/978-3-540-68585-2\\_10](http://dx.doi.org/10.1007/978-3-540-68585-2_10)
- [19] P. Viola and M. J. Jones, "Robust real-time face detection," *International Journal of Computer Vision*, vol. 57, pp. 137–154, 2004, 10.1023/B:VISI.0000013087.49260.fb. [Online]. Available: <http://dx.doi.org/10.1023/B:VISI.0000013087.49260.fb>
- [20] M. Whitty, "Liar, liar! an examination of how open, supportive and honest people are in chat rooms," *Computers in Human Behavior*, vol. 18, no. 4, pp. 343 – 352, 2002. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0747563201000590>
- [21] M. Whitty and J. Gavin, "Age/sex/location: Uncovering the social cues in the development of online relationships," *CyberPsychology Behavior*, vol. 4, no. 5, pp. 623–630, October 2001.
- [22] J. Lynch, "From fingerprints to DNA," Electronic Frontier Foundation, Tech. Rep., May 2012. [Online]. Available: <https://www.eff.org/sites/default/files/filenode/BiometricsImmigration052112.pdf>