# Reputation System for Decentralized Population Census

S. Qin[1], M. C. Silaghi[1], T. Matsui[2], M. Yokoo[3], and K. Hirayama[4]

[1]Florida Tech, [2]Nagoya Inst. of Technology, [3]Kyushu Univ., [4]Kobe Marine Univ.

**Abstract.** We describe a framework and techniques for running decentralized census processes that enable observers to independently verify governmental data. Census is a process impacting important issues such as the representation of a community in the Congress and the amount of funding that it gets from a central government. Correct census is essential for detecting vote stuffing. Reliable census data can enable user certification for addressing fake identities and Sybil attacks. Census has been historically run by governments, but citizens and NGOs need to be able to verify it. Classical census is expensive and beyond the reach of these players, hence the need for affordable citizen-driven census technology. Various citizens have different opinions as to what information should be gathered and what makes a person eligible to be accounted for in statistics. Using as inputs the official preferences of a given government should enable the verification of the data of that government. The reported work formalizes this problem and introduces a framework for reasoning about census data. An addressed challenge is how to quantify the uncertainty and the trust in the data provided by users. We report on techniques to reason and to extract census-related conclusions based on the available data. Probabilistic models with various approximations are experimented for evaluating the census results in this context.

**Keywords:** Reputation, Decentralized Census, False Identities, Sybil

## 1 Introduction

We address the problem of gathering census data using a decentralized, citizen-driven mechanism. The challenge addressed here consists in formalizing the census problem and developing algorithms for reasoning about confidence in obtained data and its implications.

Census processes have been run for thousands of years by governments as ways of estimating expected taxes and/or military power. Classically a census employs humans to talk to each resident, counting the population of a country as well as gathering certain data items about each individual, to be used in designing and justifying policy making. These processes have been run by governments and its power to alter the published figure has been identified as the main threat to stability in certain societies, as they enable large scale ballot stuffing [1]. It is therefore in the interest of citizens and non-governmental organizations (NGOs)

to verify and validate the results of the official census. For the sake of stability, it is also important for governments to increase the confidence of the citizens in official operations. Running an independent census on a large scale and using classical procedures is a complex operation and most of the interested individuals cannot afford it. Achieving census results with the quality with which governments achieve them is unlikely in the absence of similar funding. However, it can be valuable and satisfying for some observers to be able to even partly verify and corroborate official data.

*Example 1.* For example, assume an activist claims that the government miscounted his area reporting 6000 people instead of his estimation of 20000. If an independent census process trusted by this activist cannot return more than 6000 people, then the fears of the activist can be alleviated.

*Example 2.* Similarly, an activist may claim that a government has moved 20000 people in an area to change its ethnic composition. If an independent census cannot find more than a couple hundreds recent immigrants into that region, then fears may be alleviated. Alternatively, if a census identifies tens of thousands of recent immigrants, then claims can be corroborated.

The results of an independent census may complement results of an official census (as it may give new opportunities for reaching additional residents. As such, governments can improve their official data using results from independent census processes, potentially after additional verification.

*Eligibility* The counting of residents as part of a census is done differently function of the philosophical principles of the given government [9, 13, 3].

*Example 3.* In Switzerland, only citizens can vote at the federal level. However, in certain cantons (states) and municipalities, voting rights are granted to foreigners having lived there for some time, e.g., for ten years in Lausanne [10].

Should a person residing in the area for the last 4 month be counted in the census? Since eligibility varies for regions, politics and even personal beliefs, the approach in this paper is to allow different semantics of eligibility to co-exist in the data gathering process. Each end-user of the census process can compute the final statistics according to her own principles and preferences, or according to the official set of preferences (when she is testing official results).

The domain of eligibility is captured by organizations. An organization is defined by a statute that governs the way in which constituents are defined and the way in which they take decisions. Organizations can range from a club or a company, to a country.

For the management of census in large organizations we employ the concept of neighborhood. In our study neighborhoods are hierarchical, with the top of the hierarchy representing the global body of constituents. The lowest level is selected such as to define groups of constituents, preferably the largest such that each constituent in each group can verify with reasonable effort an identity

claiming to belong to that same group (what is reasonable effort may depend on the organization and could, for example, be evaluated as one day of work). Another assumptions about a neighborhood is that everybody can verify the existence of neighborhoods that are siblings in the hierarchical tree to any ancestor neighborhood of the group to which it belongs (e.g., constituents in a city can verify with reasonable effort the existence of another given locality in their county, or of another state in their country).

Constituents of these organizations verify their census continuously by witnessing on each other's eligibility in the organization. Witnessing can be seen as an open vote, but a vote based on the reputation of the voter. While witnessing for a false identity can decrease the reputation of the voter, identities witnessed by a trusted source inherit some of that reputation.

Each observer trusts itself, and this defines a rooted graph of trust on which inferences relevant to the observer can be made.

After introducing related work, in section Concepts we introduce the main definitions. Section Techniques introduces the experimented algorithms for evaluating user data. We conclude after discussing experimental results.

## 2   Background

Online White Pages directories [2] and online voter lists [5], come the closest to the task of enumerating citizens, but none of them attempts specifically to be exhaustive. People can opt to be left out of White Pages. Also, White Pages list only people having a telephone line, and may list only a subset of the inhabitants linked to a given phone line. Voter lists are typically not available freely in their entirety [5]. They contain only users that voted in previous elections. Moreover, at one moment we noticed that one can edit the information of some voters without authorization and verification. Although in this case the information is no longer collected only by companies, but also by citizens, the correctness of citizen provided information cannot be verified.

One of the main challenges of large distributed collaborations is that one user can login under as many identities as she has time and desire to register. The creation and usage of such duplicated identities is referred in literature as the Sybil attack. The term Sybil attack was first introduced by [8] in a generic distributed computing environment. In the presence of a trusted authority, the resistance to Sybil attacks is either offered by explicitly certified participation as in Microsoft's Farsite [15] or by an implicit verification. This implicit verification can be regarded as too dependent on unsafe assumptions about underlying systems, as in the Cooperative File System [7].

There is no globally trusted authority in our assumptions for the decentralized census process. That raises the problem of how to validate counterfeit identities. In [8], possible methods are classified into direct validation and indirect validation approaches. The former suggests that an entity only accepts identities that it has directly validated by some means. The latter suggests that an entity accepts identities that are vouched for by already accepted identities.

Similar techniques are used in X509 certification schemes [6]. Our approach for validation of counterfeit identities is related to this, but we bring the idea of eliciting and exploiting both positive validations and negative validations for each identity.

Census processes with validation can be successful only if people are sufficiently connected to provide enough data to the decision making process. Studies of connectivity between people have been conducted in relation to existing social networks. A kind of constituency was discussed in [11].

A reputation system maintains scores inferred from other's opinions for participants [14, 17]. Notions of valued trust are proposed in [16] and extended in [4]. The introduced values can be used to decide if an entity is sufficiently trustworthy. The values are inferred from a graph with nodes as entities and edges as the trust relations. They formalize trust relations of different types, among which are identification (ID) and trustworthiness (PR), and discuss the potential offered by networks of such relations to model known distributed authentication protocols. Bayesian models employing various probability distribution functions have been used to model behavior and its dynamics over time in Bayesian Reputation Systems [12, 17].

## 3   Framework

In this section we introduce in detail the definitions of the items involved in a decentralized census.

*Constituents* The people with right to cast votes that have a predefined weight in an organization form its constituency. Not all members of the constituency are able to input their own data (due to availability, illness, age, or lack of skills). The members that do not generate items in the virtual space are referred to as *inactive constituents*. Users that directly generate items are referred to as *active constituents*.

*Neighborhood* For ease of accounting, constituents can be organized in tree structures with nodes (called neighborhoods) corresponding to localities, cities, counties, states and countries. In this case, localities, cities, counties, states and countries form a natural hierarchy. The leaf of the tree of neighborhoods is the smallest cell of the census management, and can be configured to correspond in real life to a block, a street or an area small enough (relatively to the population density) such that members can learn and easily verify residency of their neighbors.

*Witness* Constituents in a organization can support or oppose the other constituent items' eligibility for being counted in a census. We say that they perform favorable or unfavorable witness stances for those identities. A witness stance can be associated with a set of semantic statements (as epistemological commitments, *favorable* or *unfavorable*, associated to ontological commitments from a set $\Omega$), such as:

  − existence versus nonexistence of constituent name-address pair,

- active constituent public key belongs or not to the constituent with declared name-address pair,
- favored versus disfavored version of a multiply occurring constituent (e.g., at the current residence versus an old residence, or with a correct name versus a misspelled name),
- eligibility versus ineligibility of constituent,
- correctness versus inaccuracy of details in identity,
- reliability versus sloppiness of witness.

For example, when a constituent $A$ declares constituent $B$ to be a *sloppy witness*, then $A$ believes that $B$ does not carefully verify all the constituents that it witnesses, unlike a *reliable witness*.

Constituents can also witness about the legitimacy of a neighborhood. For example, they can state that no locality called *Geneva* exists in their county, or that no street called $21^{st}$ *Street* exists in their city. Semantic statements for such witness stances can be of type:

- favored versus disfavored version of a multiply occurring neighborhood (e.g., *New York* vs *New-York city*)
- existing vs nonexistent neighborhood

The witness concept can be generalized by extending the epistemological commitments from *favorable* and *unfavorable* to something else (e.g. probabilities). The current article is focused on the simpler case with *favorable* and *unfavorable* and the generalizations are left for other studies.

## 4  Census Process Concepts

Now let us introduce concepts involved in the decentralized census processes.

### 4.1  Citizen Interactions

A citizen-driven census requires participation of individual citizens for actions such as *residence declaration* and *witnessing*. As residence declarations, each individual voluntarily provides census data not only about herself but also about her neighbors. The neighborhood where a citizen resides is part of its identity details. A voting process, called witnessing, is used to help verify the census data. The verification can be done both by neighbors, and by volunteers who gather data about the inhabitants of the given area.

*Witness Graph* A graph defined by the witness relations between constituents can be generated in the following way:

- A node is generated for each constituent.
- A directed edge from node $A$ to node $B$ is generated for each semantic statement that $A$ witnesses for $B$.

- Each edge has a color (from a set $\Omega$), given by the type of statement that generated it (ontological commitment).
- An edge has weight 1 if generated for a favorable stance and weight 0 if generated for an unfavorable stance (epistemological commitment).

Inactive nodes are sinks for this graph. This graph can be used to reason about the eligibility of the declared identities and implicitly about the census.

*Distributed Census Problem* The Distributed Census Problem (DCP) for an observer $\Gamma$ can be formalized as a tuple $\langle \mathcal{N}_\mathcal{S}, \mathcal{I}, \mathcal{R}, \mathcal{W}, \mathcal{M}_\mathcal{S}, \Gamma \rangle$, where:

- $\mathcal{N}_\mathcal{S}$ is the set of neighborhoods $\mathcal{N}_\mathcal{S} = \{1, ..., d\}$,
- $\mathcal{I}$ is the set of person identities,
- $\mathcal{R}$ is the set of residence declarations (constituent items)
- $\mathcal{W}$ is the set of witness stances
- $\mathcal{M}_\mathcal{S}$ is a model of the relation between the ground truth $I^*$ and $\mathcal{N}_\mathcal{S}$, $\mathcal{I}$, $\mathcal{R}$ and $\mathcal{W}$, as believed by the observer $\Gamma$ (e.g., a certain belief network)

$\mathcal{I}^*$ (the ground truth), each having an identity from the set $\mathcal{I}$. The problem is to approximate the $\mathcal{I}^*$ that best explains $\mathcal{N}_\mathcal{S}$, $\mathcal{I}$, $\mathcal{R}$ and $\mathcal{W}$ based on the model $\mathcal{M}_\mathcal{S}$.

## 5 Techniques

Here we present the techniques used to address the challenge of inferring a count of the constituency given a witness graph.

*Eligibility* Although anyone can participate in the census process of an organization, not everyone is eligible to be counted in the census. In an organization, which is the context of this study, the definition of eligibility is a function of the constituent. When the eligibility for a constituent is based on a subjective view, the census result is relevant only to the user (or users) sharing this view. Hence, we define the eligibility as a probabilistic function of several parameters:

- Someone's interpretation of the witness graph, $\mathcal{M}_\mathcal{S}$
- Someone's own definition of the eligibility, $\Gamma(\mathcal{O})$

**Definition 1.** *The reference user is the user $\Gamma$ who currently computes the census.*

**Definition 2 (Censable and $\Psi$).** *A constituent item $C$ is* `censable` *for an organization if it is eligible and new (never counted elsewhere). The $\Gamma$'s confidence value in whether $C$ is* `censable` *is denoted $\Psi(C)$.*

**Definition 3 (Witness Reliability and $\Phi$).** *A constituent item $C$ is a* reliable witness *if $\Gamma$ trusts all the witness stances that $C$ issues as she trusts her own. $\Gamma$ may not fully trust the stances of another constituent $C$, but only with a confidence value $\Phi(C)$.*

Based on the DCP parameters, one can infer a value $\Psi$ for the confidence that observer $\Gamma$ can have on whether a given constituent item $C$ identifies a `censable` user, and a value $\Phi$ for its confidence on whether $C$ is *witness reliable.*

*Remark 1 (Decision Criteria 1).* One approach to compute a census is to declare that an identity is eligible (to be counted in the census) from the point of view of the reference user $\Gamma$ if the value of $\Psi$ surpasses a threshold $t$ where $t$ is defined by the $\Gamma$.

*Remark 2 (Decision Criteria 2).* Another approach is to sum the values $\Psi$ for all constituents (once normalized in the interval $[0,1]$).

We introduce the notation $\Phi$ to denote the quality *reliable witness* when used as superscript or subscript with one of the notations above (e.g., $SC^{\Phi}(C)$). Similarly we use $\Psi$ to denote the `censable` quality when used as superscript or subscript in these notations (e.g., $SC^{\Psi}(C)$).

### 5.1   Probabilistic Models

We will now discuss models based on Bayesian Networks. In this model, random variables are used to represent the `censable` property of each constituent, the *reliable witness* property, and the witnessing stances between each pair of constituents for each quality. All these random variables are Boolean. For each pair of constituents $A$ and $B$ we get the random variables and Bayesian Network in Figure 1. Note that each pair of constituents requires the introduction of $2|\Omega|$ random variables for $|\Omega|$ qualities. With the two considered qualities in Figure 1, $\Phi$ and $\Psi$, a constituent (e.g., $A$) is associated with two hidden random variables: `censable` ($CS^A$) and *reliable_witness* ($RW^A$). Each pair of constituents items (e.g., A and B) is associated with four evidence (grayed) random variables: A witnesses for B being a reliable witness ($W^{AB\Phi}$), A witnesses for B being `censable` ($W^{AB\Psi}$), B witnesses for A being a reliable witness ($W^{BA\Phi}$), B witnesses for A being `censable` ($W^{BA\Psi}$).

While conditional probability tables can be trained from real data once large amount of such data is available, sample conditional probability tables built manually for variables of type $W^{BA\Phi}$ and $W^{BA\Psi}$ are shown in Table 1.

| $CS^A$ | $RW^B$ | $P(W^{BA\Psi})$ |
|---|---|---|
| $t$ | $t$ | 0.9 |
| $t$ | $f$ | 0.5 |
| $f$ | $t$ | 0.1 |
| $f$ | $f$ | 0.5 |

| $RW^A$ | $RW^B$ | $P(W^{BA\Phi})$ |
|---|---|---|
| $t$ | $t$ | 0.9 |
| $t$ | $f$ | 0.5 |
| $f$ | $t$ | 0.3 |
| $f$ | $f$ | 0.5 |

| $CS^A$ | $RW^A$ | $P(W^{AA\Psi})$ |
|---|---|---|
| $t$ | $t$ | 0.99 |
| $t$ | $f$ | 0.5 |
| $f$ | $t$ | 0.1 |
| $f$ | $f$ | 0.5 |

| $RW^A$ | $P(W^{AA\Phi})$ |
|---|---|
| $t$ | 0.99 |
| $f$ | 0.5 |

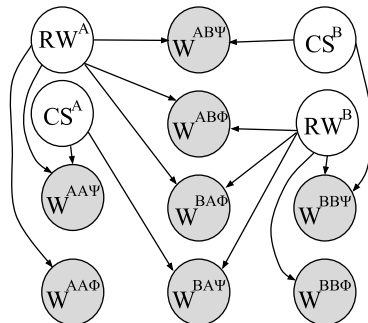| $P(RW)$ | $P(CS)$ |
|---|---|
| 0.5 | 0.5 |

**Table 1.** Transition, sensor and prior CPT

**Fig. 1.** A Bayesian Network for two users

**Theorem 1.** *The number of random variables (all Boolean) in a Bayesian Network modeling a DCP is linear in the size of the input.*

*Proof.* Given $n$ constituent items $C_1, ..., C_n$, we get at most $n^2|\Omega|$ random variables:

- $n|\Omega|$ hidden variables modeling the real qualities of each constituent item, and
- $n(n-1)|\Omega|$ modeling the evidence variables about all $|\Omega|$ possible witness stances between each of the $n(n-1)$ possible directed pairs of constituents.

Note that we do not need to model with random variable the nonexistent witness stances. Therefore the actual network size is linear in the size of the input, being proportional to the number $w$ of input semantic statements in witness stances $(w + n)|\Omega|$.

For average sized networks one can perform queries of values for the random variables $CS^{C_i}$, modeling $\Psi(C_i)$ of $i^{th}$ constituent item, using techniques such as Markov Chain Monte Carlo (MCMC).

## 6    Experiments

To evaluate the power of the studied DCP models to represent users reasoning about census, as well as to resist various attacks, we perform two sets of preliminary experiments. One of them is based on a set of volunteers and the second is based on a larger simulated data.

In the experiments based on volunteers we asked 10 people living within an area of a few square kilometers to register themselves as active constituents and to also register others 10 friends as inactive constituents of a regional organization. Each of these volunteers had the opportunity to witness for the other constituent items that they knew. We also introduced 2 obviously wrong constituents at an address that most participants knew to not exist. A snapshot

of the interactions between constituents is shown in Figure 2 where the thick edges represent favorable witness stances, the thin edges represents unfavorable witness stances, the nodes represents constituents.
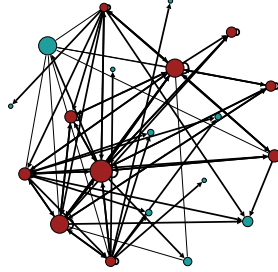


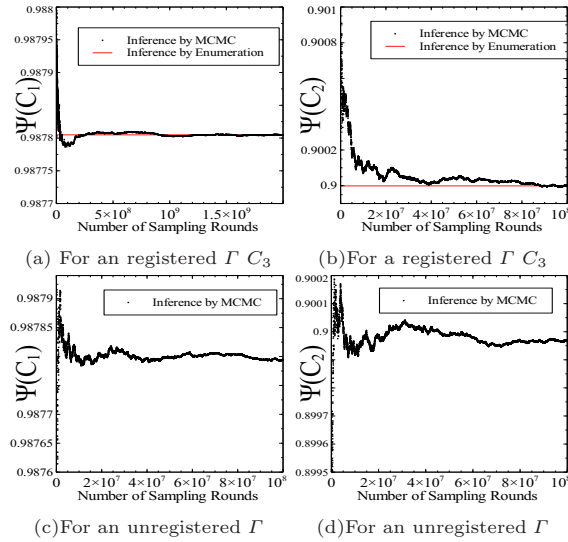**Fig. 2.** Visualization of constituents and the witness relation between them



**Fig. 3.** Inference of $\Psi(C_1)$ and $\Psi(C_2)$

With the Bayesian network described in Figure 1 and CPT specified by Table 1, we use MCMC to perform queries of values for the random variable $CS^{C_i}$ which models the $\Psi(C_i)$ of the $i^{th}$ users. Convergence for a few users ($C_1$ and $C_2$) is shown in Figures 3 (a), (b), (c), (d). The reference (red line) is computed with exact inference by enumeration (which at this problem size was possible within a day). For an unregistered $\Gamma$, the exact inference by enumeration is expected to take 33 days on a computer and is not shown here.

We have performed extensive experiments with simulated data, targeting the evaluation of the robustness to attacks that such a system can provide.

Most of them will not be described here for lack of space. The result of a set of 5 experiments estimating the impact of the percentage $k$ of honest active constituents(HACs) in the global population on the `censable` properties of the constituents is shown in Figure 4. The true positive rate (TPR) gives the percentage of correctly counted constituents out of the total number of eligible constituents. The false positive rate (FPR) gives the percentage of wrongly counted constituents with respect to the total number of eligible constituents. A robust census process has a high TPR and a low FPR. The semantic statements for witness stances are only about the eligibility quality. In this experiment, a constituent item $(n, a)$ is eligible if there is someone whose name is $n$ and lives at address $a$. We simulate attackers that declare a number of ineligible constituent items and perform favorable witness stances for a percentage of $h$ ineligible constituent items in their leaf neighborhoods. We assume that witness stances represent all semantic statements. The plotted points are for $h$ of 100%, 93.75%, 87.5%, 75%, 50%, 25%, 12.5%, 6.25%, 0% respectively. The studied values of $k$ are (0.9, 0.8, 0.7, 0.5 and 0.3) respectively.
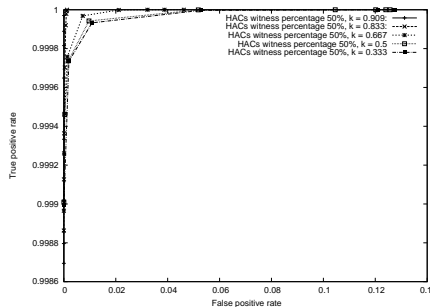


**Fig. 4.** The effects of the percentage of honest active constituent ($k$ value) in global population with curves defined by varying $h$.

There are several common parameters for the plotted curves. The total number of eligible constituents is 9300000. The *HACs witness percentage* is 50% out of their leaf neighborhoods. *HACs witness percentage* is the percentage of constituent items that an HAC witnesses honestly (A favorable witness stance is performed if an item is eligible and an unfavorable witness stance is performed if an item is ineligible). The number of *attackers* is 300000. The number of ineligible constituent items declared by each *attackers* is 4.

Since we see in Figure 4 that the curve with parameter $k$=0.9 is higher than the curve with parameter $k$= 0.8, the curve with parameter $k$=0.8 is higher than the curve with parameter $k$=0.7, the curve with parameter $k$=0.7 is higher than the curve with parameter $k$=0.5 and the curve with parameter $k$=0.5 is higher
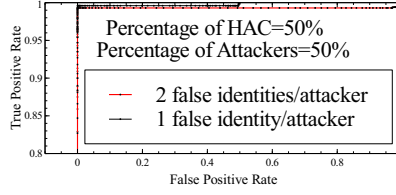
**Fig. 5.** ROC with 1000 MCMC rounds, 1000 constituents and no neighborhoods

than the curve with parameter $k=0.3$, we conclude that the $k$ value will affect the robustness of the system positively. That is, the bigger the $k$ value is, the more accurate the system will be.

Figure 5 illustrates the ROC curve for an unstructured experiment (without neighborhoods) based on 1000 real constituents and 1000 MCMC rounds. Each attacker creates $n$ fake identities into the global population ($n \in \{1, 2\}$) and has a favorable witness stance on each of the fake identities. Percentages of HACs and attackers out of all constituents are both 50%. Each HAC has a correct witness on 1% of all identities (real and false). Constituent $C$ is counted here when $\Psi(C)$ is greater than a threshold $t$ (varying between 0 and 1). The ROC curves reveals a low false positives rate and high true positives rate, indicating a robust system of this configuration of parameters. The curve for more Sybils is slightly below the other one, which is consistent with graceful degradation.

## 7   Adversarial Learning

The system is implemented in DirectDemocracyP2P, where it is used for detecting false identities. Each user has a full copy of all the data and can independantly compute the census with the parameters of her choice. An adversary can try to learn the parameters used by most users and adapt its behavior to manipulate the outcome. For example, one of the behaviors of an adversary would be to gain trust from peers with high connectivity. Somebody's trust could be easier obtained by replicating that users's votes.

## 8   Conclusions

We have addressed the problem of formalizing and solving the decentralized population census problem (DCP). While population census is an important process with large implications in the distribution of public funds and security of elections from vote stuffing, it is currently an expensive process outside the reach of external verifiers and was identified as a threat to stability in certain regions. To enable a decentralized citizen-driven census, we investigate a set of concepts such as: organization, constituent, neighborhood and witnessing. The organization is a set of rules (constitution) the specify mechanisms to define eligibility of constituents. For large organizations, the constituency is organized in a tree of

neighborhoods to help with census organization. Constituents can witness (vote) on each other's qualities, such as: eligibility and witnessing reliability.

We have also proposed and analyzed theoretically a probabilistic model based on Bayesian Networks that can be used to address the problem in a principled way. Experiments with large simulated data show that robustness to attackers is possible when there exists a reasonable kernel of honest active constituents.

## References

1. O. Analytica. Politics of census-taking threatens stability. `http://www.oxan.com/display.aspx?ItemID=DB154173`, 2009.
2. ATT Interactive Inc. `http://www.yellowpages.com/whitepages`, 2010.
3. L. Beckman. Citizenship and voting rights: Should resident aliens vote? *Citizenship studies*, 10(2):153–165, 2006.
4. T. Beth, M. Borcherding, and B. Klein. Valuation of trust in open networks. *Computer Security-ESORICS 94*, pages 1–18, 1994.
5. Complex Systems Inc. `http://www.floridavoterfile.com`, 2010.
6. D. Cooper, S. Santesson, S. Farrell, R. Housley, and W. Polk. Internet x.509 public key infrastructure certificate and certificate revocation list (crl) profile. *Network Working Group Request for Comments*, 5280, May 2008.
7. F. Dabek, M. Kaashoek, D. Karger, R. Morris, and I. Stoica. Wide-area cooperative storage with CFS. In *Symposium on OS principles*, pages 202–215. ACM, 2001.
8. J. Douceur. The sybil attack. *Peer-to-peer Systems*, pages 251–260, 2002.
9. D. Earnest. Noncitizen voting rights: A survey of an emerging democratic norm. In *Annual convention of the American Political Science Assoc.*, pages 28–31, 2003.
10. D. C. Earnest. Voting rights for resident aliens: A comparison of 25 democracies. In *Annual Meeting of the Northeast Political Science Association and the International Studies Association-Northeast*, 2003.
11. J. Fedoruk, A. ; Denzinger. A general framework for multi-agent search with individual and global goals: Stakeholder search. *International Transactions on Systems Science and Applications (ITSSA)*, 1(4):357–362, 2006.
12. A. Jøsang and W. Quattrociocchi. Advanced features in bayesian reputation systems. In *Trust, Privacy and Security in Digital Business*, pages 105–114. 2009.
13. D. Owen. Resident aliens, non-resident citizens and voting rights: towards a pluralist theory of transnational political equality and modes of political belonging. In *Citizenship Acquisition and National Belonging: Migration, Membership and the Liberal Democratic State*, pages 52–73. November 2009.
14. P. Resnick, K. Kuwabara, R. Zeckhauser, and E. Friedman. Reputation systems. *Communications of the ACM*, 43(12):45–48, 2000.
15. Thawte. Web of trust. `http://www.thawte.com/secure-email/web-of-trust-wot/`, 2013.
16. R. Yahalom, B. Klein, and T. Beth. Trust relationships in secure systems-a distributed authentication perspective. In *Symposium on Research in Security and Privacy*, pages 150–164. IEEE, 1993.
17. M. U. Zeinab Noorian. The state of the art in trust and reputation systems: A framework for comparison. *JTAER*, 5(2):97–117, 2010.