# Beyond De-Identification Record Falsification to Disarm Expropriated Data-Sets

Roger Clarke

# Beyond De-Identification
# Record Falsification to Disarm Expropriated Data-Sets

ROGER CLARKE

**Abstract** The wild enthusiasm for big data and open data has brought with it the assumptions that the utility of data-sets is what matters, and that privacy interests are to be sacrificed for the greater good. As a result, techniques have been devised to reduce the identifiability of expropriated data-records, on the assumption that privacy is to be compromised to the extent necessary. This paper argues for and adopts data privacy as the objective, and treats data utility for secondary purposes as the constraint. The inadequacies of both the concept and the implementation of de-identification are underlined. Synthetic data and Known Irreversible Record Falsification (KIRF) are identified as the appropriate techniques to protect against harm arising from expropriated data-sets.

**Keywords:** • Big Data • Data Analytics • Privacy • Re-Identification • Record Falsification •

CORRESPONDENCE ADDRESS: Roger Clarke, Xamax Consultancy Pty Ltd, Australian National University, Research School of Computer Science, Sydney, Australia, e-mail: Roger.Clarke@xamax.com.au.

## 1       Introduction

During the early decades of administrative computing, roughly 1950-1980, personal data was collected for a purpose, used for that purpose, and confined to 'silos'.  Only in particular circumstances did it escape beyond its original context, and become subjected to 're-purposing', and combination with data from other sources.  From the 1970s onwards, however, there was growth in the financial services sector's sharing of data about consumers' creditworthiness (Furletti 2002), and in data matching by government agencies (Clarke 1994c).

Over the last few decades, these, initially exceptional, secondary uses of personal data have changed from a dribble to a haemorrhage, supported by advances in the technical capabilities necessary to handle large volumes of data.  The notions of 'data warehousing' (Inmon 1992) and 'data mining' (Fayyad et al. 1996) emerged.  After early disappointments, these ideas have recently resurged with the new marketing tags of 'big data', 'open data', 'data analytics' and 'data science'.

During the early decades of data protection law, the fundamental principle was that use and disclosure beyond the original purpose of collection had to be based on consent or authority of law (OECD Use Limitation Principle, OECD 1980). The protection that this Principle was meant to afford has since been torn asunder by exemptions, exceptions and long lists of legal authorisations.

In this paper, the short-form term 'expropriation' is adopted to refer to the kinds of secondary use that are common with big data / open data.  These enthusiastic movements are based on the application of data for purposes beyond the aims for which it was collected and is authorised by the individual to be used.
In the public sector, governments around the world appear to have been inspired by the openness of Danish agencies' databases to secondary uses (e.g. Thygesen et al. 2011).  A practice is becoming more widespread in which all personal data that is gathered by government agencies, in many cases under compulsion, is regarded as the property of the State and "a strategic national resource that holds considerable value" (AG 2015).  The whole of government is treated as a monolith – thereby breaching the 'data silo' protection mechanism.

In Australia, for example, the Australian Institute of Health and Welfare (AIHW) has pillaged data-sources across the healthcare sectors at national and state and

territory levels, and made rich sub-sets available to large numbers of researchers. Further, the Australian Bureau of Statistics (ABS) has a multi-agency data integration program (MADIP) in train with a range of 'partner' agencies. This extracts data gathered for specific administrative purposes and enables its analysis for a wide range of purposes. A great many other such projects are being conducted and proposed under the big data and open data mantras. Justifications for the abuses of data in government and in health-related research emphasise collectivism and de-value individualism.

The private sector is piggybacking on the 'open data' notion (e.g. Deloitte 2012). Corporations are encouraged by governments to treat data about individuals as an exploitable asset, irrespective of its origins, sensitivity and re-identifiability. Assertions of business value, and that such activities are good for the economy, are treated as being of greater importance than human values.

In both sectors, proponents and practitioners make the assumption that such projects are capable of delivering significant benefits, even though tha data has been wrenched far beyond its original context, has been merged with other data with little attention paid to incompatibilities of meaning and quality, and has been analysed for purposes very different from those for which it was collected. Limited attention is paid to data quality audit and even less to testing of the inferences drawn from such data-collections against real-world patterns (Clarke 2016b). Considerable scepticism is necessary about the real effectiveness and social value of these activities.

The doubts extend beyond the activities' justification to the negative impacts on the individuals whose data is expropriated. Proponents of big data do not object to replacing identifiers with pseudonyms; but they do not welcome comprehensive privacy protection: "it is difficult to ensure the dataset does not allow subsequent re-identification of individuals, but ... it is also difficult to de-identify datasets without introducing bias into those sets that can lead to spurious results" (Angiuli et al. 2015). Significantly for the argument pursued in this paper, the position adopted by big-data proponents is that the interests of the individuals to whom the expropriated data relates are secondary, and that such procedures as are applied to reduce the risk of harm to individuals' privacy must be at limited cost to its utility for organisations.

Examples of the claim for supremacy of the data-utility value over the privacy value abound. For example, "we underline the necessity of an appropriate research exemption from consent for the use of sensitive personal data in medical research ..." (Mostert et al. 2016, emphasis added) More generally, "We develop a method that allows the release of [individually identifiable microdata] while minimizing information loss and, at the same time, providing a degree of preventive protection to the data subjects" (Garfinkel et al. 2007, p.23, emphasis added).

The theme for the Bled conference in 2019 is 'Humanising Technology for a Sustainable Society'. This paper addresses that theme by proposing a switch back from the asserted supremacy of data utility to recognition of the primacy of the human right of privacy. It is not argued that data utility should be ignored. The proposition is that, when preparing personal data for disclosure and use beyond its original context, the appropriate value to adopt as the objective is privacy protection. The retention of such utility as the data may have for other purposes is not the objective. It remains, however, an important factor to be considered in the choice among alternative ways of ensuring that harm is precluded from arising from re-identification of the data.

The paper commences with a summary of privacy concerns arising from the expropriation of personal data and its use and disclosure for purposes far beyond the context within which it was collected. The notions of identification, de-identification and re-identification are outlined, and conventional techniques described. This builds on a long series of prior research projects by the author. De-identification is shown to be a seriously inadequate privacy-protection measure. Two appropriate approaches are identified: synthetic data, and Known Irreversible Record Falsification (KIRF).

The paper's contributions are the review of de-identification measures from the perspective of the affected individuals rather than of the expropriating parties, and the specification of falsification as a necessary criterion for plundered data-sets.

## 2 The Vital Role of Data Privacy

Privacy is a pivotal value, reflected in a dozen Articles of the International Covenant on Civil and Political Rights (ICCPR 1966). It underpins many of the rights that are vital constituents of freedom. Fuller discussion is in Clarke (2014c). Philosophical analyses of privacy are often based on such precepts as human dignity, integrity, individual autonomy and self-determination, and commonly slide into conflicts between the moral and legal notions of 'rights'. Adding to the confusion, legal rights vary significantly across jurisdictions. A practical working definition is as follows (Morison 1973, Clarke 1997):

**Privacy** *is the interest that individuals have in sustaining 'personal space', free from interference by other people and organisations*

The diversity of contexts within which privacy concerns arise is addressed by typologies that identify dimensions or types of privacy (Clarke 1997, Finn et al. 2013, Koops et al. 2016). The dimensions of privacy of personal data and of personal communications are directly relevant to the present topic. The term 'information privacy' is commonly used to encompass both data at rest and on the move, and is usefully defined as follows:

**Information privacy** *is the interest that individuals have in controlling, or at least significantly influencing, the handling of data about themselves.*

Protection of information privacy is not only important in its own right. It also provides crucial underpinning for protections of the other three dimensions: privacy of personal behaviour, of personal experience, and of the physical person.

Abuse of the privacy interest results in significant harm to human values. Within communities, psychological harm and negative impacts on social cohesion are associated with loss of control over one's life and image, loss of respect, and devaluation of the individual. Reputational harm inflicted by the disclosure of data about stigmatised behaviours, whether of the individual or of family-members, reduces the pool of people prepared to stand for political office and hence weakens the polity. Profiling, and use of data-collections to discover behaviour-patterns and generate suspicion, lay the foundation for the repression of behaviours that powerful organisations regard as undesirable. This

undermines the exposure of wasteful, corrupt and otherwise illegal activities, and reduces the scope for creativity in economic, social, cultural and political contexts. At any given time, a proportion of the population is at risk of being identified and located by a person or organisation that wishes to take revenge against them or exact retribution from them, excite mortal fear in them, or eliminate them.

Behavioural privacy is harmed not only from unjustified collection, use and disclosure of personal data, but also from the knowledge or suspicion that individuals may be watched, that data may be collected, and that their activities may be monitored. This has a 'chilling effect' on group behaviour, whereby intentional acts by one party have a strong deterrent effect on important, positive behaviours of some other party (Schauer 1978). This results in stultification of social and political speech. A society in which non-conformist, inventive and innovative behaviour are stifled risks becoming static and lacking in cultural, economic and scientific change (Kim 2004).

Data sensitivity is relative. Firstly, it depends on the personal values of the individual concerned, which are influenced by such factors as their cultural context, ethnicity, lingual background, family circumstances, wealth, and political roles. Secondly, it depends on the individual's circumstances at any particular point in time, which affects what they want to hide, such as family history, prior misdemeanours, interests, attitudes, life-style, assets, liabilities, or details of their family or family life.

Various aspects of privacy are important, in particular circumstances, for a substantial proportion of the population. Some categories of individual are more highly vulnerable than others. For the large numbers of people who at any given time fall within the many categories of 'persons-at-risk', it is essential to guard against the disclosure of a great deal of data, much of it seemingly innocuous (GFW 2011, Clarke 2014a).

To assist in assessment of the effectiveness of safeguards against harm arising from data expropriation, Table 1 presents a small suite of test-cases that are sufficiently diverse to capture some of the richness of human needs.

**Table 1: Six 'Persons-at-Risk' Test-Cases**

- **People with outlier, non-conformist or 'deviant' personal profiles**
  Key Data:  characteristics of interest to service-providers
  Key Risks: denial of service
  e.g. genetic or medical conditions resulting in discrimination by health insurers, low 'social credit' scores resulting in denial of access to transport

- **Negotiators of corporate mergers and acquisitions**
  Key Data:  information-sources, locations, meeting-partners
  Key Risks: breach of corporations law and stock exchange listing rules

- **Candidates for political office**
  Key Data:   associations with stigmas such as psychiatric treatment
  Key Risks:   unelectability, reduction in the pool of candidates

- **Whistleblowers and media sources**
  Key Data:   identity
  Key Risks:   retribution, drying-up of informers, unchecked corruption

- **Victims of domestic violence**
  Key Data:  location
  Key Risks:   physical harm

- **Police informants and protected witnesses**
  Key Data:   pseudonym and/or location
  Key Risks:   physical harm, loss of witness, loss of future witnesses

This paper's purpose is to switch the focus away from the asserted utility of big and open data for secondary purposes, and back towards the human value of privacy.  However, there are further aspects of the conference theme of 'Humanising Technology for a Sustainable Society' that are negatively affected by the prevalence of data expropriation.

Many organisations' operations depend on access to personal data, and on the quality of that data. For an analysis of data quality aspects in big data contexts, see Clarke (2016b). The goodwill of the individuals concerned is very important not only to data access and data quality, but also to the cost incurred in assuring data quality. Extraneous uses of personal data cause a significant decrease in trust by individuals in the organisations that they deal with. The result is that they are much less willing to disclose and much more likely to hide and to obscure data, and much less willing to disclose honestly, and much more likely to disclose selectively, inconsistently, vaguely, inaccurately, misleadingly, imaginatively or fraudulently. There is a great deal of scope for obfuscation and falsification of data (Schneier 2015a, 2015b, Bösch et al. 2016, Clarke 2016a). Widespread exercise of these techniques will have serious negative consequences for the quality of data held by organisations.

Expropriation of data results in the data on which analyses are based bearing a less reliable relationship to the real-world phenomena that they nominally represent. This leads to the inferences that are drawn by medical, criminological and social research in the interests of the public, and by marketing activities in the interests of corporations, being at best misled and misleading, and their use being harmful rather than helpful. This particular form of dehumanising technology, rather than contributing to the sustainability of society, undermines it.

This section has presented the reasons why privacy is a vital human value. The proposal that privacy is the primary objective and data-utility the constraint is therefore of far more than merely academic interest, and is a social and economic need. The following section outlines the relevant aspects of identification, and the conventional mechanisms that have been applied to expropriated data in order to achieve what designers have portrayed as being 'anonymisation' of the data.

## 3    De- and Re-Identification

This section outlines the notions of identity, nymity and identification, drawing on Clarke (1994b, 2010). It presents the conditions that need to be fulfilled in order that de-identification can be achieved, and re-identification precluded. It then provides an overview of techniques applied to expropriated personal data.

## 3.1    Concepts

An **entity** is a real-world thing. Rather than artefacts such as tradeable items and mobile phones, this paper is concerned with human beings. An **identity** is an entity of virtual rather than physical form. Each person may present many identities to different people and organisations, and in different contexts, typically associated with roles such as consumer, student, employee, parent and volunteer. During recent decades, organisations have co-opted the term 'identity' to refer to something that they create and that exists in machine-readable storage. Better terms exist to describe that notion, in particular **digital persona** (Clarke 1994a, 2014b). In this paper, the term 'identity' is used only to refer to presentations of an entity, not to digital personae.

The notion of **'nymity'** is concerned with identities that are not associated with an entity. In the case of **anonymity**, the identity cannot be associated with any particular entity, whether from the data itself, or by combining it with other data. On the other hand, **pseudonymity** applies where the identity is not obviously associated with any particular entity, but association may be possible if legal, organisational and technical constraints are overcome (Clarke 1999).

An **identifier** is a data-item or set of data-items that represent attributes that can reliably distinguish an identity from others in the same category. Commonly, a human identity is identified by name (including context-dependent names such as 'Sally' or 'Herbert' at a service-counter or in a call-centre), or by an identifying code that has been assigned by an organisation (such as an employee- or customer-number).

**Identification** is the process whereby a transaction or a stored data-record is associated with a particular identity. This is achieved by acquiring an identifier, or assigning one, such as a person's name or an identifying code.

**De-identification** notionally refers to a process whereby a transaction or a stored data-record becomes no longer associable with a particular identity. However, it is in practice subject to a number of interpretations, outlined in Table 2.

**Table 2: Alternative Interpretations of 'De-Identification'**

1. **The removal of data-items** that are designed to, or are known to, facilitate the association of a record with a real-world identity. This interpretation is the one most commonly apparent in the literature. It satisfies a necessary condition, but falls a long way short of being sufficient

2. **Further adaption and/or 'perturbation' of the data-set in order to address additional association risks**. These are discovered by analysis of the data and its various contexts in order to achieve understanding of the many other ways in which at least some proportion of the records may remain associable with a particular real-world identity. This interpretation is sometimes apparent in the literature

3. **Further processing of the data-set to address the risk of physical or virtual merger, linkage or comparison of that data-set with other data-sets**. This interpretation is seldom apparent in the literature

4. **Demonstration of the reliability of de-identification**, by showing that the records in the data-set cannot be associated with the real-world identity to whom they originally applied. This interpretation is seldom apparent in the literature

De-identification of a data-set is very likely to result in at least some degree of compromise to the data-set's utility for secondary purposes. In Culnane et al. (2017), it is argued that "decreasing the precision of the data, or perturbing it statistically, makes re-identification gradually harder at a substantial cost to utility". It remains an open question as to whether, under what circumstances, and to what extent, the objectives of the two sets of stakeholders can be reconciled. For early examinations of the **trade-off between de-identification and the utility of the data-set**, see Duncan et al. (2001), Brickell & Shmatikov (2009) and Friedman & Schuster (2010). The perception of compromise to data utility appears to be an important reason why the more powerful de-identification techniques in Table 2 are seldom actually applied, or at least not with the enormous care necessary to achieve significant privacy-protection.

In many circumstances, de-identified records are subject to 're-identification', that is to say the re-discovery or inference of an association between a record and a real-world identity, despite prior attempts to de-identify them. This is possible because de-identification is extremely difficult for all but the simplest and least interesting data-sets. It is particularly easy with rich data-sets, such as those whose records contain many data-items, or whose data-items contain unusual values.

Further, a great many of the data-sets that are lifted out of their original context and re-purposed are subsequently merged or linked with other data-sets. This gives rise to two further phenomena, which together greatly increase the risk of inappropriate matches and inappropriate inferences (Clarke 2018):

- combined data-sets generally offer even more opportunities for re-identification than do single-source data-sets; and
- combined data-sets are far more likely than single-source data-sets to lead to faulty inferences being drawn. This is because:
    - the quality of the data in each of the data-sets is often not high and hence comparisons of data-content may be unreliable;
    - the meanings of the data-items in each of the data-sets are often unclear or ambiguous;
    - the definitions of the data-items in each of the data-sets may be inconsistent or otherwise incompatible; and
    - where data scrubbing activities have been undertaken, before and/or after combination of the data-sets, the process(es) of addressing some problems inevitably also create new problems.

The notion of re-identification has attracted considerable attention, particularly since it was demonstrated that "87% ... of the population in the United States had reported characteristics that likely made them unique based only on {5-digit ZIP, gender, date of birth}. About half of the U.S. population ... are likely to be uniquely identified by only {place, gender, date of birth}, where place is basically the city, town, or municipality in which the person resides" (Sweeney 2000).

Narayanan & Shmatikov (2008) presented a general de-anonymization algorithm which they claimed requires "very little background knowledge (as few as 5-10 attributes in our case study). Our de-anonymization algorithm is robust to imprecision of the adversary's background knowledge and to sanitization or perturbation that may have been applied to the data prior to release. It works even if only a subset of the original dataset has been published" (p.2). For fuller discussion of re-identification, see Ohm (2010) and Slee (2011).

## 3.2    Longstanding De-Identification Techniques

In response to objections to the expropriation of personal data, proponents argue that the records are 'anonymised', can no longer be associated with the individual concerned, can therefore do that individual no harm, and hence the individual should not be concerned about the re-use or disclosure of the data. In order to deliver what they claim to be 'anonymised data', expropriating organisations have applied a variety of techniques.

From the 'data mining' phase (indicatively 1980-2005), a literature exists on **'privacy-preserving data mining'** (PPDM – Denning 1980, Sweeney 1996, Agrawal & Srikant 2000). For a literature review, see Brynielsson et al. (2013). PPDM involves suppressing all identifiers and other data-items likely to enable re-identification ('quasi-identifiers'), and editing and/or statistically randomising (or 'perturbing') the contents of data-items whose content may assist re-identification (e.g. because of unusual values). The declared purpose is to preserve the overall statistical features of the data, while achieving a lower probability of revealing private information.

During the later 'big data' phase (since c. 2010), guidance on forms of data manipulation that are suitable for practical application is provided in particular by UKICO (2012), but see also DHHS (2012). In Slee (2011), a simple set of four categories is suggested: **replacement, suppression, generalisation and perturbation**. Accessible summaries of the challenges and some of the risks involved in the de-identification process are in Garfinkel (2015) and Polonetsky et al. (2016).

The regulatory regime applying to US health records (HIPAA) specifies two alternative approaches for de-identification: The Expert Determination Method and the Safe Harbor method (which is effectively a simplified 'fool's guide').

However, "neither method promises a foolproof method of de-identification with zero risk of re-identification. Instead, the methods are intended to be practical approaches to allow de-identified healthcare information to be created and shared with a low risk of re-identification" (Garfinkel 2015, p. 22).

In D'Acquisto et al. (2015) pp.27-37, it is noted that "most data releasers today (e.g. national statistical offices) tend to adopt the *utility-first approach*, because delivering useful data is their *raison d'etre*" (p.29, emphasis added). A further indication of the strong commitment to data utility is that, although "in Germany, any organizational data accessible to external researchers is required to be de facto anonymized", the bar is set very low, because all that is required is that "the *effort* that is necessary to identify a single unit in the data set is *higher than the actual benefit* the potential intruder would gain by the identification" (Bleninger et al., 2010, emphasis added). This formulation ignores the critical issues that (1) the breach causes harm to the affected individual, and (2) the harm to the affected individual may be far greater than the benefit to the breacher.

Similarly, the de-identification decision-making framework in O'Keefe et al. (2017) remains committed to the utility-first approach, because it applies the threshold test of "when data is *sufficiently de-identified* given [the organisation's] data situation" (p.2, emphasis added).

A further indicator of the inadequacy of the approaches adopted is that 're-identification risk' is regarded as being merely "the percentage of de-identified records that can be re-identified" (Garfinkel 2015, p. 38). If privacy rather than utility is adopted as the objective, then 're-identification risk' is seen to be a much more complex construct, because every breach has to be evaluated according to the potential harm it gives rise to – which can be severe in the case of a wide range of categories of persons-at-risk.

Garfinkel's conclusion was that, **"after more than a decade of research, there is comparatively little known about the underlying science of de-identification"** (2015, p.39). Given the complexities involved in both the problems and the techniques, it is far from clear that any practical solutions will ever emerge that satisfy the privacy-first rather than the utility-first criterion.

## 3.3      Re-Identification Techniques

The application of de-identification techniques naturally stimulated responses: "it seems that new techniques for de-identifying data have been met with equally innovative attempts at re-identification" (Hardy & Maurushat 2017, p.32). For analyses of techniques for re-identification, see Sweeney (2002), Acquisti & Gross (2009) and Ohm (2010).

In relation to one critical area of concern, the re-identifiability of location and tracking data, Song et al. (2014) showed that "human mobility traces are highly identifiable with only a few spatio-temporal points" (p.19). Further, de De Montjoye et al. (2015) found that credit card records with "four spatiotemporal points are enough to uniquely reidentify 90% of individuals ... [and] knowing the price of a transaction increases the risk of reidentification by 22%" (p. 536). Culnane et al. (2017) and Teague et al. (2017) described successful re-identification of patients in a de-identified open health dataset.

In contesting the De Montjoye et al. findings, Sanchez et al. (2016) provided a complex analysis, concluding that "sound anonymization methodologies exist to produce useful anonymized data that can be safely shared ...". It is inconceivable that the intellectual effort brought to bear by those authors in defending disclosure would or even could ever be applied to the continual, high-volume disclosures that are part-and-parcel of the data expropriation economy: "De-identification is not an exact science and ... you will not be able to avoid the need for complex judgement calls about when data is sufficiently de-identified given your [organisation's] data situation" (O'Keefe et al. 2017, p.2). The practical conclusion is that, at least where privacy is prioritised over data-utility:

*Sound anonymization methodologies are so complex and onerous that they cannot be relied upon to produce useful anonymized data that can be safely shared*

The re-identification process is easier where:

1.   the data-set contains a large number of data-items;
2.   there are unique values within individual data-items; and/or
3.   there are unique combinations of values across multiple data-items.

A further important consideration is the availability of multiple data-sets that are capable of being compared, which gives rise to greater richness in a combined or merged data-set. An important factor in successful de-identification activities has been the widespread availability of large data-sets, such as electoral rolls, subscription lists, profiles on social networking sites, and the wide range of data broker offerings. In short, **a great many expropriated data-sets satisfy the conditions for easy re-identification of a material proportion of the records they contain.**

## 3.4    Recent De-Identification Techniques

The D'Acquisto monograph describes more privacy-protective techniques that have been proposed by academics – although most of them appear to be encountering difficulty in escaping the laboratory. The monograph refers to the alternative techniques as 'privacy-first anonymisation', but use of that term is not justified. The formulation is still utility-as-objective and privacy-as-constraint: "a parameter ... guarantees an upper bound on the re-identification disclosure risk and perhaps also on the attribute disclosure risk". Further, even in academic experimentation, the privacy-protectiveness has been set low, due to "*parameter choices relaxing privacy* in order for reasonable utility to be attainable" (p.29, emphasis added).

The D'Acquisto et al. summary of the 'privacy models' underlying these techniques is as follows: "A first family includes k-anonymity and its extensions taking care of attribute disclosure, like p-sensitive k-anonymity, l-diversity, t-closeness, (n,t)-closeness, and others. The second family is built around the notion of differential privacy, along with some variants like crowd-blending privacy or BlowFish" (D'Acquisto et al. 2015, p.30).

The k-anonymity proposition is a framework for quantifying the amount of manipulation required of quasi-identifiers in order to achieve a given level of privacy (Sweeney 2002). A data-set satisfies k-anonymity iff each sequence of values in any quasi-identifier appears with at least k occurrences. Bigger k is better. The technique addresses only some of the threats, and has been subjected to many variants and extensions in an endeavour to address further threats.

Differential privacy is a set of mathematical techniques that reduces the risk of disclosure by adding non-deterministic noise to the results of mathematical operations before the results are reported. An algorithm is differentially private if the probability of a given output is only marginally affected if one record is removed from the dataset (Dwork 2006, 2008).

In both cases, "The goal is to keep the data 'truthful' and thus provide good utility for data-mining applications, while *achieving less than perfect privacy*" (Brickell & Shmatikov 2009, p.8, emphasis added). Further, the techniques depend on assumptions about the data, about other data that may be available, the attacker, the attacker's motivations, and the nature of the attack. Some of the claims made for the techniques have been debunked (e.g. Narayanan & Shmatikov 2010, Zang & Bolot 2011, Narayanan & Felten 2016, Zook et al. 2017, Ashgar & Kaafar 2019), and a range of statistical attacks is feasible (O'Keefe & Chipperfield 2013, pp. 441-451). **All k-anonymity and differential privacy techniques provide very limited protection.**

Even if these highly complex techniques did prove to satisfy the privacy-first criterion, the excitement that they have given rise to in some academic circles has not been matched in the real world of data expropriation, and it appears unlikely that they ever would be. None of the techniques, nor even combinations of multiples of them, actually achieve the objective of privacy-protection – not least because their aim is the retention of the data's utility. The highest standard achieved within the data-utility-first tradition, even in the more advanced, but seldom implemented forms, might be reasonably described as 'mostly de-identified' or 'moderately perturbed'.

The data-utility-first approach, and the de-identification techniques that it has spawned, cannot deliver adequate privacy protection. The expropriation of personal data gives rise to harm to people generally, and is particularly threatening to person-at-risk such as the small suite of test-cases in Table 1. Addressing their needs requires another approach entirely.

# 4       Privacy-First Disarming of Expropriated Data-Sets

This section considers ways in which privacy can be prioritised, but, within that constraint, such utility as is feasible can be rescued from data-sets. Although it is unusual for researchers to treat privacy as the objective and economic benefits as the constraint, it is not unknown. For example, in Li & Sarkar (2007), "The proposed method attempts to preserve the statistical properties of the data based on privacy protection parameters specified by the organization" (p.254). Privacy is thereby defined as the objective, and the statistical value of the data the constraint ("attempts to preserve").

In another approach, Jändel (2014) describes a process for analysing the risk of re-identification, and determining whether a given threshold ("the largest acceptable de-anonymisation probability for the attack scenario") is exceeded. When the safety of victims of domestic violence and protected witnesses is taken into account, that threshold has to be formulated at the level of impossibility of discovery of the person's identity and/or location. It is therefore reasonable to treat Jändel's extreme case as being privacy-protective.

In order to provide adequate protection against privacy breaches arising from expropriated data-sets even after de-identification, two approaches are possible:

1.   Avoid the risks, by not using empirical data, but instead generating synthetic data
2.   Prevent the risks arising, by ensuring that, even where individual records are re-identified, the data is unusable because it has been falsified in ways the specifics of which are unknowable, and which are irreversible

The remainder of this section considers those two approaches.

## 4.1     Synthetic Data

The most obvious way in which privacy can be protected is by not expropriating data, and hence avoiding use and disclosure for secondary purposes. This need not deny the extraction of utility from the data. Under a variety of circumstances, it is feasible to create 'synthetic data' that does not disclose data that relates to any individual, but that has "characteristics that are similar to real-world data

[with] frequency and error distributions of values [that] follow real-world distributions, and dependencies between attributes [that are] modelled accurately" (Christen & Pudjijono 2009. p.507).

This has been argued by some to be an effective solution to the problem: "empirically, it is difficult to find a database table on which sanitization permits both privacy and utility. Any incremental utility gained by non-trivial sanitization (as opposed to simply removing quasi-identifiers or sensitive attributes) is more than offset by a decrease in privacy, measured as the adversarial sensitive attribute disclosure. It is possible, however, to construct an artificial database, for which sanitization provides both complete utility and complete privacy, even for the strongest definition of privacy ..." (Brickell & Shmatikov 2009, p.7).

To date, there appears to have been very little take-up of this approach. As abuses of personal data, and harm arising from them, become increasingly apparent to the public, the assumed power of national statistical and other government agencies and large corporations may be shaken, and the generation of synthetic data may become much more attractive.

## 4.2    Empirical Data, De-Fanged

In this case, the proposition is that no data-set can be expropriated beyond its original context unless it has been first rendered valueless for any purpose relating to the administration of relationships between organisations and particular individuals. One way of achieving this is to convert all record-level data that was once empirical – in the sense of being drawn from and reflecting attributes of real-world phenomena – into synthetic data that represents a plausible phenomenon, but not a real one.

The underlying data-set is of course not affected, and remains in the hands of the organisation that manages it. The underlying data-set is the appropriate basis for administering the relationships between organisations and particular individuals;  whereas expropriated data-sets are not.

The process must also be irreversible, at the level of each individual data record.

Further, **the fact of processing (as distinct from the details), and the standards achieved:**

- **must be known by organisations that do or may gain access to the expropriated data-sets.** This ensures that they are aware that the record-level data, whether or not it can be associated with any particular person, is unusable for any purpose related to the individual; and

- **must be known by affected individuals, and by advocacy organisations for their interests.** This ensures confidence in the process, and avoids motivating people to obfuscate or falsify data about themselves

Combining these properties, this mechanism is usefully described as **Known Irreversible Record Falsification (KIRF)**.

The possibility exists that the characteristics of some data-sets, or of some records within them, may resist falsification to the point of unusability. In that case, the records in question are unsuitable for expropriation, and no empirical derivative of them may be disclosed. If those records constitute a sufficient proportion of the data-set as a whole, then the data-set as a whole cannot be disclosed.

Examples of data-sets that may contain records that are too rich to be effectively falsified include the combination of psychological and social data with stigmatised medical conditions, and data about undercover operatives in national security and law enforcement contexts. (This of course does not necessarily preclude the use of statistical distributions derived from such data-sets as a basis for generating synthetic data that has comparable overall characteristics).

A corollary of the privacy-first approach is that the utility of the data-set is a constraint, not an objective. This might seem to rob the expropriated data-set of a great deal of value. Intuitively, it would appear unlikely that any single process could achieve both the standard of 'irreversibly falsified records' and preservation of the original data-set's overall statistical features. On the other hand, for any given use to which the expropriated data-set is to be put, different falsification processes could be applied, in order to produce a data-set that

preserves the particular statistical features that are critical for that particular analysis.

In most circumstances, it would appear likely that changes can be made to data in order to satisfy the criteria, while sustaining at least a moderate level of utility for particular purposes. This is an empirical question that cannot be determined in the abstract, but requires detailed analysis in each specific context of data-set and purpose.

A less stringent approach could be considered, whereby the 'every record' requirement is relaxed, in favour of 'enough records'. However, because many records are not falsified, the data-set's utility for making decisions about individuals is not undermined and hence adversaries are motivated to conduct attacks. Individuals whose records are not falsified are subject to compromise. This is a serious matter, because inevitably some of them would be among the categories of persons-at-risk. The inadequacy extends further, however, because the interests of all individuals are compromised. Records that have been falsified are also likely to be used to generate inferences – and, due to the falsification steps, the inferences that are drawn are unreliable, and potentially harmful.

The less stringent arrangement would fail to curb the eagerness of organisations to exploit the expropriated data-set, and would fail to earn the trust of the affected individuals. Even if the application of a particular record's content to a particular individual were to be precluded by law, the scope for unregulated abuse of the provision is too high. The Known Irreversible Falsification criterion needs to be applied to all records, not merely to some or even most records.

## 5 Towards an Evaluation Process for the Privacy-First Approach

The purpose of this paper has been to argue for a privacy-first approach to the preparation of data-sets for expropriation to secondary purposes, and to develop an operational definition of what that involves. This section provides some preliminary suggestions as to the steps necessary to apply the principles, operationalise the process, and assess its effectiveness.

The term 'privacy-first' is of recent origin, and its first use in D'Acquisto et al. (2015, p.29) was in any case a false start. The sense in which it is proposed in this paper is so far outside the present mainstream as to be arguably deviant. Searches for existing literature on data perturbation undertaken to satisfy the requirement of falsification have not located a literature on the topic, or even individual instances that adopt the approach. Further, in the absence of theoretical discussions, it is not likely that exemplars and testbeds can be readily found.

On the other hand, some prior work is very likely to have relevance, in the sense of being capable of adaptation to the privacy-first criterion. A simple example of this would be a model in which a parameterisation mechanism enables the privacy weighting to be set at 1, but that nonetheless delivers non-zero utility, or at least information or insights. An approach to generating action in this field would be to expose these ideas in workshops that focus on de-identification and re-identification topics.

It is also feasible for projects to be undertaken that commence with existing guidelines on data perturbation, apply the Known Irreversible Record Falsification (KIRF) principle, and test the results by considering the 6 test-cases in Table 1. Initial projects might use data-sets of convenience. More serious studies would then be needed on mainstream, rich data-sets, such as those in the Census, social data and health care fields that are commonly subjected to expropriation.

Once the point has been reached that multiple approaches have been specified that satisfy the requirement, further rounds of research are needed in order to establish principles and practical guidance in relation to the retention of maximal utility, while still satisfying the requirement of known irreversible falsification for all individual records.

## 6      Conclusion

This paper's purposes have been:

- to abandon the utility-first approach;
- to adopt privacy as the objective and relegate data-utility to the level of a constraint;
- to argue for data-expropriation beyond its original context to be contingent on the prior application of techniques that fulfil that requirement;  and
- to identify and articulate specific ways in which this can be done.

The analysis of alternative criteria for achieving privacy-first disarming of data-sets identified two contenders.  The first possibility is the use only of synthetic data.  This avoids the disclosure of any personal data, by creating and disclosing data whose distribution has usefully close approximations to the original data, but without any scope for disclosure of any personal data relating to any actual identity.

The second possibility applies the Known Irreversible Record Falsification (KIRF) criterion, in order to achieve similar properties in a released data-set to those of synthetic data.  This achieves privacy protection by ensuring that all records are unusable for any purpose that relates to any specific individual.  KIRF will, however, have impacts on the utility of data-sets.  These impacts may be modest, but will often be significant, and in some circumstances will render the data-set unusable for data analytics purposes.

An implication of this conclusion is that research into de-identification processes needs to shift away from the approaches adopted over the last 15 years, such as k-anonymity and differential privacy, which prioritise utility at the expense of privacy. Instead, **the need is for a focus on ways to minimise the harm to the utility of data-sets, given that every record has to be falsified in such a manner that it is unusable for determinations about individuals, and is known to be so.**

If data-expropriating organisations fail to switch their approach in this way, it will be increasingly apparent to the public that their personal data is being expropriated and exploited by organisations without meaningful regard for either the rights of individuals or the harm that may arise from re-identification. As one research team in the re-identification area put it, "The ... government holds vast quantities of information about [individuals]. It is not really 'government data'. It is data about people, entrusted to the government's care" (Culnane et al. 2017).

A proportion of the population will neither know nor care. A further proportion will know, and care, but feel themselves to be technically incapable and/or powerless to do anything about it, sullenly accept the situation, and trust organisations as little as possible. The remainder will take action in order to deny the use of their data. Many techniques are already been demonstrated whereby individuals can resist abuse of their data, and moderate numbers of tools for obfuscation and falsification are available for deployment.

Over the last 50 years, organisations' data-gathering techniques have migrated from manual capture by employees to a combination of manual capture by the individuals to whom the data relates and automated capture as a byproduct of transactional activities. There is increasing incidence of autonomous creation of data by equipment that individuals are not aware are monitoring their behaviour. Obfuscation and falsification are easiest in relation to the long-standing forms of data capture. There are interesting challenges aplenty in devising ways to avoid, subvert and defeat byproduct and autonomous data capture. The expertise of many capable individuals will be attracted to the endeavour.

If this scenario unfolds, the quality of data that is in the expropriated collections will diminish below its present mediocre level. This will have serious implications for the validity, and for the business and policy value, of inferences drawn from data analytics activities. The benefits to economies and societies arising from this scenario will be significantly less than what would be achieved if instead the privacy-first approaches advocated above are adopted. This paper thereby contributes to the aim of humanising technologies for sustainable society.

## References

Acquisti A. & Gross R. (2009) `Predicting Social Security Numbers from Public Data' Proc. National Academy of Science 106, 27 (2009) 10975-10980

AG (2015) 'Public Data Policy Statement' Australian Government, December 2015

Agrawal R. & Srikant R. (2000) 'Privacy-preserving data mining' ACM Sigmod Record, 2000

Aggarwal C.C. & Philip S.Y. (2008) 'Privacy-Preserving Data Mining: Model and Algorithms' Springer, 2008

Angiuli O., Blitzstein J. & Waldo J. (2015) 'How to De-Identify Your Data' Communications of the ACM 58, 12 (December 2015) 48-55

Asghar H.J. & Kaafar D. (2019) 'Averaging Attacks on Bounded Perturbation Algorithms' arxiv.org, February 2019

Bleninger P., Drechsler J. & Ronning G. (2010) 'Remote data access and the risk of disclosure from linear regression: An empirical study' In 'Privacy in Statistical Databases', Vol. 6344 of Lect Notes Comp Sci, Eds. J. Domingo-Ferrer & E. Magkos. Berlin Heidelberg: Springer; pp.220–233

Bösch C., Erb B., Kargl F., Kopp H. & Pfattheicher S. (2016) 'Tales from the Dark Side: Privacy Dark Strategies and Privacy Dark Patterns' Proc. Privacy Enhancing Technologies 4 (2016) 237-254

Brickell J. & Shmatikov V. (2009) 'The Cost of Privacy: Destruction of Data-Mining Utility in Anonymized Data Publishing' Proc. KDD'08, August 2008

Brynielsson J., Johansson F. & Jändel M. (2013) 'Privacy-preserving data mining: A literature review' Swedish Defence Research Agency, February 2013

Christen P. & Pudjijono A. (2009) 'Accurate synthetic generation of realistic personal information' Proc. PAKDD, Springer LNAI, vol. 5476, pp. 507–514

Clarke R. (1994a) 'The Digital Persona and its Application to Data Surveillance' The Information Society 10,2 (June 1994) 77-92, PrePrint at http://www.rogerclarke.com/DV/DigPersona.html

Clarke R. (1994b) 'Human Identification in Information Systems: Management Challenges and Public Policy Issues' Information Technology & People 7,4 (December 1994) 6-37, PrePrint at http://www.rogerclarke.com/DV/HumanID.html

Clarke R. (1994c) 'Dataveillance by governments: The technique of computer matching' Information Technology & People 7, 2 (December 1994), 46-85, PrePrint at http://rogerclarke.com/DV/MatchIntro.html

Clarke R. (1997) 'Introduction to Dataveillance and Information Privacy, and Definitions of Terms' Xamax Consultancy Pty Ltd, August 1997, revisions to July 2016, at http://www.rogerclarke.com/DV/Intro.html

Clarke R. (1999) 'Identified, Anonymous and Pseudonymous Transactions: The Spectrum of Choice' Proc. User Identification & Privacy Protection Conference, Stockholm, 14-15 June 1999, PrePrint at /DV/UIPP99.html http://www.rogerclarke.com/DV/UIPP99.html

Clarke R. (2008) 'Dissidentity: The Political Dimension of Identity and Privacy' Identity in the Information Society 1, 1 (December, 2008) 221-228, PrePrint at http://www.rogerclarke.com/DV/Dissidentity.html

Clarke R. (2010) 'A Sufficiently Rich Model of (Id)entity, Authentication and Authorisation' Proc. IDIS 2009 - The 2nd Multidisciplinary Workshop on Identity

in the Information Society, version of February 2010, at http://www.rogerclarke.com/ID/IdModel-1002.html

Clarke R. (2014a) 'Key Factors in the Limited Adoption of End-User PETs' Xamax Consultancy Pty Ltd, April 2014, at http://www.rogerclarke.com/DV/UPETs-1405.html#PU

Clarke R. (2014b) 'Promise Unfulfilled: The Digital Persona Concept, Two Decades Later' Information Technology & People 27, 2 (Jun 2014) 182 - 207, PrePrint at http://www.rogerclarke.com/ID/DP12.html

Clarke R. (2014c) 'Privacy and Free Speech' Invited Presentation at the Australian Human Rights Commission Symposium on Free Speech, Sydney, 7 August 2014, Xamax Consultancy Pty Ltd, August 2014, at http://www.rogerclarke.com/DV/PFS-1408.html

Clarke R. (2016a) 'A Framework for Analysing Technology's Negative and Positive Impacts on Freedom and Privacy' Datenschutz und Datensicherheit 40, 1 (January 2016) 79-83, PrePrint at http://www.rogerclarke.com/DV/Biel15-DuD.html

Clarke R. (2016b) 'Quality Assurance for Security Applications of Big Data' Proc. EISIC'16, Uppsala, 17-19 August 2016, PrePrint at http://www.rogerclarke.com/EC/BDQAS.html

Clarke R. (2018) 'Guidelines for the Responsible Application of Data Analytics' Computer Law & Security Review 34, 3 (May-Jun 2018) 467- 476, PrePrint at http://www.rogerclarke.com/EC/GDA.html

Culnane C., Rubinstein B.I.P. & Teague V. (2016) 'Understanding the maths is crucial for protecting privacy' Pursuit, 29 September 2016

Culnane C., Rubinstein B.I.P. & Teague V. (2017) 'Health Data in an Open World' arXiv, December 2017

D'Acquisto G., Domingo-Ferrer J., Kikiras P., Torra V., de Montjoye Y.-A. & Bourka A. (2015) 'Privacy by design in big data' ENISA, December 2015

Deloitte (2012) 'Open data: Driving growth, ingenuity and innovation' Deloitte, London, 2012

de Montjoye Y.-A., Hidalgo C.A., Verleysen M. & Blondel V.D. (2013) 'Unique in the Crowd: The privacy bounds of human mobility' Sci. Rep. 3, 1376 (2013)

De Montjoye Y.-A., Radaelli L., Singh V.K. & Pentland A. (2015) 'Unique in the Shopping Mall: On the Reidentifiability of Credit Card Metadata' Science 347, 6221 (29 January 2015) 536–539

Denning D.E. (1980) 'Secure statistical databases with random sample queries' ACM TODS 5, 3 (Sep 1980) 291- 315

DHHS (2012) 'Guidance Regarding Methods for De-identification of Protected Health Information in Accordance with the Health Insurance Portability and Accountability Act (HIPAA) Privacy Rule' Department of Health & Human Services, November 2012

Duncan G.T., Keller-McNulty S.A. & Stokes S. L. (2001) 'Disclosure Risk vs Data Utility: The R-U Confidentiality Map' Technical Report LA-UR-01-6428, Los Alamos National Laboratory

Dwork, C. (2006) 'Differential Privacy' In: Bugliesi, M., Preneel, B., Sassone, V., Wegener, I. (eds.) ICALP 2006. LNCS, vol. 4052, pp. 1–12. Springer, Heidelberg (2006)

Dwork C. (2008) 'Differential Privacy: A Survey of Results' in M. Agrawal et al. (Eds.):

TAMC 2008, LNCS 4978, pp. 1–19, 2008

Fayyad U., Piatetsky-Shapiro G. & Smyth P. (1996) 'From Data Mining to Knowledge Discovery in Databases ' AI Magazine 17, 3 (1996) 37-54

Finn R. L., Wright D. & Friedewald M. (2013) 'Seven types of privacy' in Gutwith S. et al. (eds) 'European Data Protection: Coming of Age' , Springer, 2013, pp. 3-32

Friedman A. & Schuster A. (2010) 'Data Mining with Differential Privacy' Proc. KDD'10, July 25–28, 2010

Furletti M. (2002) ' An Overview and History of Credit Reporting' Federal Reserve Bankl of Philadelphia, June 2002

Garfinkel R., Gopal R. & Thompson S. (2007) 'Releasing Individually Identifiable Microdata with Privacy Protection Against Stochastic Threat: An Application to Health Information' Information Systems Research 18, 1 (Mar 2007) 23-41,121-122

Garfinkel S.L. (2015) ' De-Identification of Personal Information' NISTIR 8053, National Institute of Standards and Technology, October 2015, at https://nvlpubs.nist.gov/nistpubs/ir/2015/nist.ir.8053.pdf

GFW (2011) 'Who is harmed by a "Real Names" policy?' Geek Feminism Wiki, undated, apparently of 2011, at http://geekfeminism.wikia.com/wiki/Who_is_harmed_by_a_%22Real_Names%22_policy%3F

Hardy K. & Maurushat A. (2017) 'Opening up government data for Big Data analysis and public benefit' Computer Law & Security Review 33 (2017) 30-37

ICCPR (1996) 'International Covenant on Civil and Political Rights' United Nations, 1966

Inmon B. (1992) 'Building the Data Warehouse' Wiley, 1992

Jändel M. (2014) 'Decision support for releasing anonymised data' Computers & Security 46 (2014) 48-61

Kim M.C. (2004) 'Surveillance technology, Privacy and Social Control' International Sociology 19, 2 (2004) 193-213

Koops B.J., Newell B.C., Timan T., Korvanek I., Chokrevski T. & Gali M. (2016) 'A Typology of Privacy' University of Pennsylvania Journal of International Law 38, 2 (2016)

Li X. & Sarkar S. (2007) 'Privacy Protection in Data Mining: A Perturbation Approach for Categorical Data' Information Systems Research 17, 3 (Sep 2006) 254-270,321-322

Morison W.L. (1973) 'Report on the Law of Privacy' Government Printer, Sydney, 1973

Mostert M., Bredenoord A.L., Biesaart M.C.I.H. & van Delden J.J.M. (2016) 'Big Data in medical research and EU data protection law: challenges to the consent or anonymise approach' Eur J Hum Genet 24, 7 (July 2016) 956–960

Narayanan A. & Felten E.W. (2016) 'No silver bullet: De-identification still doesn't work' In Data protection on the move 2016 (pp. 357–385), Springer Netherlands

Narayanan A. & Shmatikov V. (2008) 'Robust De-anonymization of Large Sparse Datasets' Proc. IEEE Symposium on Security and Privacy, pp. 111-125

Narayanan A. & Shmatikov V. (2010) 'Myths and fallacies of personally identifiable information' Commun. ACM 53,6 (June 2010) 24-26

OECD (1980) 'OECD Guidelines on the Protection of Privacy and Transborder Flows of Personal Data' OECD, Paris, September 1980

Ohm P. (2010) 'Broken Promises of Privacy: Responding to the Surprising Failure of

Anonymization' 57 UCLA Law Review 1701 (2010) 1701-1711

O'Keefe C.M. & Chipperfield J.O. (2013) 'A summary of attack methods and confidentiality protection measures for fully automated remote analysis systems' International Statistical Review 81, 3 (December 2013) 426–455

O'Keefe C.M., Otorepec S., Elliot M., Mackey E. & O'Hara K. (2017) 'The De-Identification Decision-Making Framework' CSIRO Reports EP173122 and EP175702

Polonetsky J., Tene O. & Finch K. (2016) 'Shades of Gray: Seeing the Full Spectrum of Practical Data De-Identification' Santa Clara Law Review 56 (2016) 593

Sanchez D., Martinez S. & Domingo-Ferrer J. (2016) 'Comment on 'unique in the shopping mall: on the reidentifiability of credit card metadata'' Science, 351, 6279 (18 March 2016)

Schauer F. (1978) 'Fear, Risk and the First Amendment: Unraveling the Chilling Effect' Boston University Law Review 58 (1978) 685-732

Schneier B. (2015a) 'Data and Goliath: The Hidden Battles to Collect Your Data and Control Your World' Norton, March 2015

Schneier B. (2015b) 'How to mess with surveillance' Slate, 2 March 2015

Slee T. (2011) 'Data Anonymization and Re-identification: Some Basics Of Data Privacy: Why Personally Identifiable Information is irrelevant' Whimsley, September 2011

Song Y., Dahlmeier D. & Bressan S. (2014) 'Not So Unique in the Crowd: a Simple and Effective Algorithm forAnonymizing Location Data' Proc. Workshop on privacy-preserving IR, pp, 19-24, 2014

Sweeney L. (1996) 'Replacing personally-identifying information in medical records, the Scrub system' Journal of the American Medical Informatics Association (1996) 333-337

Sweeney L. (2000) 'Simple Demographics Often Identify People Uniquely' Data Privacy Working Paper 3, Carnegie Mellon University, 2000

Sweeney L. (2002) 'k-anonymity: a model for protecting privacy' International Journal on Uncertainty, Fuzziness and Knowledge-based Systems 10, 5 (2002) 557-570

Teague V., Culnane C. & Rubinstein B.I.P. (2017) ' The simple process of re-identifying patients in public health records' Pursuit, 18 December 2016

Thygesen L.C., Daasnes C., Thaulow I. & Bronnum-Hansen H. (2011) 'Introduction to Danish (nationwide) registers on health and social issues: Structure, access, legislation, and archiving' Scandinavian Journal of Public Health 39, 7 (2011) 12–16

UKICO (2012) 'Anonymisation: managing data protection risk: code of practice' Information Commissioners Office, UK, November 2012, esp. pp. 21-27, App. 2 pp. 51-53, and Annex 3 pp. 80-102 - the last section by Yang M., Sassone V. & O'Hara K., Uni. of Southampton

Zang H. & Bolot J. (2011) 'Anonymization of Location Data Does Not Work: A Large-Scale Measurement Study' Proc. MobiCom'11, September 2011

Zook M. et al. (2017) 'Ten simple rules for responsible big data research' Editorial, PLOS Computational Biology, 30 March 2017