

# Runners, repeaters, strangers and aliens: operationalising efficient output disclosure control

**Kyle Alves**

*University of the West of England (UWE), Bristol*

**Felix Ritchie**

*University of the West of England (UWE), Bristol*

# Runners, repeaters, strangers and aliens: operationalising efficient output disclosure control

Kyle Alves<sup>1</sup> and Felix Ritchie<sup>2</sup>

<sup>1</sup>Senior Lecturer of Operations and Information Systems, University of the West of England.

[Kyle.alves@uwe.ac.uk](mailto:Kyle.alves@uwe.ac.uk)

<sup>2</sup>Professor of Applied Economics, Bristol Business School, University of the West of England.

[Felix.ritchie@uwe.ac.uk](mailto:Felix.ritchie@uwe.ac.uk)

## Abstract:

Statistical agencies and other government bodies are increasingly using secure remote research facilities to provide access to sensitive data for research as an efficient way to increase productivity. Such facilities depend on human intervention to ensure that the research outputs do not breach statistical disclosure control (SDC) rules.

Output SDC can be either principles-based, rules-based, or ad hoc. Principles-based is often seen as the gold standard when viewed in statistical terms, as it improves both confidentiality protection and utility of outputs. However, some agencies are concerned that the operational requirements are too onerous for practical implementation, despite the evidence to the contrary.

This paper argues that the choice of output checking procedure should be seen through an operational lens, rather than a statistical one. We take a standard model of operations management which focuses on understanding the nature of inputs, and apply it to the problem of output checking. We demonstrate that the principles-based approach addresses user and agency requirements more effectively than either the rules-based or ad hoc approaches, and in a way which encourages user buy-in to the process. We also demonstrate how the principles-based approach can be aligned with the statistical and staffing needs of the agency.

JEL classification: C00, C80

Keywords: confidentiality, statistical disclosure control, operations management, output checking

## 1. Introduction

In the last two decades, one of the key growth areas in official statistics has been the availability of confidential data for research user by academics, private sector analysts and government departments. On the demand side, users want increasing granularity in the data to address more specific policy issues. On the supply side, government data holders are under pressure to leverage their investment in data collection by maximising data use across a range of stakeholders.

Much of this data is confidential and personal, such as health or tax data. Traditionally, the privacy of respondents was managed by reducing the detail in the data, either to a level at which the data could be distributed without restriction (public use files, or PUFs), or with more detail left in the data but access limited to licensed users (scientific use files, or SUFs).

As data use has grown, so have concerns about whether the confidentiality protection is adequate. The new risks include (Statistics Authority, 2018) the re-identification possibilities of social media, the third-party holding of confidential data implied by the growth in administrative data as a source, and massive computing power with the ability to re-identify source data through brute force methods. There have already been examples of anonymization methods which were adequate some years ago that no longer meet acceptable standards.

There appear to be five solutions to this, according to observed practice. The first is to reduce detail further; this risks making the data valueless. A second is to tighten up on the contracts for SUFs, but this does not solve the problem of PUF re-identification risk; it also assumes that there is a linear relationship between strict licensing conditions and user behaviour, for which there is no strong evidence. A third option is to replace genuine data with synthetic data, but users are often uncomfortable about basing analysis on imputed data. The fourth solution is 'query servers', systems which allow simple queries on the data with confidentiality checks applied to outputs. Table servers, producing simple cross-tabulations and counts, are becoming widespread and effective at meeting many users' needs for dynamic tabulations. More complex query servers offering a much wider range of analysis are now being developed, such as Statistics Norway's elegant system at [www.microdata.no](http://www.microdata.no).

However, for detailed analysis researchers need access to the full microdata, and so the fifth solution is to allow this in an environment under the control of the data holder – the research data centre (RDC). The great success story of this century for official statistics has been the use of virtual RDCs (vRDCs), where thin client technology has allowed data holders to provide the security of a physically restricted environment whilst allowing users to access the environment from more convenient locations. Most European countries have at least one facility operated by the National Statistics Institute (NSI) or a data archive, as do the US, Canada, Mexico, South Africa, Japan, Australia and New Zealand. In the UK alone there are six general-purpose vRDCs offering the microdata underlying official statistics to a variety of users in government and academia.

These so-called secure use files (SecUFs) address the issue of confidentiality at the point of access, but create a new risk of confidentiality breach through publication (Lowthian and Ritchie, 2017). If the data has some identification risk (as in both SUFs and SecUFs) then it is possible that a published output might reveal some confidential information. This risk is higher for SecUFs as the data is much more detailed. All RDCs therefore operate a system of output-checking before publication (output statistical disclosure control, or OSDC) to manage this risk.

There are two approaches to managing output-checking for conformance to regulation: 'rules-based' and 'principles-based' (Ritchie and Elliott, 2015). The former sets strict rules for releasing output and applies simple yes/no criteria; the latter uses flexible rules-of-thumb and creates an environment for negotiation between researcher and output-checker. Because rules-based is very limiting in research environments, our experience is that most organisations claiming to be rules-based operate a 'rules-based but sometimes...' system allowing for ad hoc relaxation of rules.

This can be viewed as a problem of risk management: which system reduces risk most? However, most data holders focus upon the operational question of efficiency: which approach uses resources most effectively? In particular, a rules-based system can, in theory be run automatically, or by humans with little statistical training; the principles-based solution requires input by humans who are able to discuss technical matters with researchers. Prima facie, principles-based seems a more costly and laborious solution, and the management literature has long established this to be the case for bespoke production (Chase, 1981). However, as ONS (2019) points out, the principles-based solution was designed *specifically* to reduce resource cost while also reducing risk, and the little evidence that is available tends to support this.

There are two reasons for the misperception of the principles-based model. First, data holders are often unfamiliar with the activities of research users of data, and so view them through the lens of their own outputs; these are typically tabulations which have strict rules applied for comparability across time and alternative breakdowns. Second, data holders' experience of OSDC is usually limited to the statistical literature, which focuses on arbitrary 'intruders' (e.g. Hundepool et al., 2010) applying mechanical procedures to breach confidentiality. Together, these factors encourage an over-simplistic view of the research environment which drives data-holders' perceptions of risk and benefits.

To illuminate this debate, we introduce a model familiar to operations management literature: that of 'runners-repeaters-strangers-aliens' (RRSA) (Parnaby, 1988; Aitken et al., 2003). This model segments inputs of demand from customers (in this case, the requests from researchers for data cleared for publication) and uses the different characteristics of those segments to develop optimal operational responses. Using this framework, we contrast how the rules-based and principles-based approaches address the different challenges posed by real research environments. It is then straightforward to demonstrate how the "one-size-fits-all" rules-based model achieves neither operational efficiency nor effective risk reduction. Similarly, we can also analyse why the "rules-based-but..." approach fails to achieve the operational advantages of the full principles-based approach.

The next section summarises the literature on the topic; this is negligible on the rules-based versus principles-based argument, but there is an extensive management literature on the RRSA model. In section three we develop the output-checking problem, and in section four we show how the RRSA model can be applied to this procedure. Section five discusses empirical cost assessments. Section six concludes.

While acknowledging that many government departments produce data for re-use by researchers in academia and government, for clarity in this article we assume that the data has been collected and made available by a national statistical institute (NSI).

## 2. Literature review

### Output checking

Output statistical disclosure control (OSDC) is a relatively new field. Until recently, the SDC literature focused almost exclusively on two problems: anonymization of microdata, and protection of tabular outputs; see for example Willenborg and de Waal (1996), or the *Privacy in Statistical Databases* biennial conference publication. Since the development of RDCs in the early 2000s, a small number of papers began to appear considering particular outputs such as regressions (Reiter, 2003; Reznik, 2004; Reznik and Riggs, 2005; Ritchie, 2006; Corscadden et al., 2006, for example) as well as general guidelines for users of RDCs (Corscadden et al., 2006).

The concept of SDC for outputs generally, and research environments in particular, was introduced in Ritchie (2007) and followed up by the concept of 'safe outputs' (Ritchie, 2008), usually referred to now as 'safe statistics' (Ritchie, 2014) or 'high/low review statistics' (ONS, 2019). Brandt et al. (2010) used these and operational practices to produce the first widely-available general purpose guide to OSDC. This was included as a chapter in Hundepool et al. (2010)'s broadly successful attempt to provide an overview of state-of-the-art techniques across the field of SDC.

Brandt et al. (2010) has been widely adopted by RDC managers as the only general guide for practitioners. It has been updated since (Bond et al., 2015) but, with the exception of Ritchie (2019) few of its precepts have undergone critical challenge. It is the main source for most subsequent publications (e.g. Eurostat 2015; Statistics NZ, 2015; O’Keefe et al., 2015).

Part of the reason for the unquestioning acceptance is the report’s attitude to the clearance process. Brandt et al. (2010) contains the first practitioner guide to both principles-based OSDC (PBOSDC), rules-based OSDC (RBOSDC), and the practical differences in implementation. Brandt et al. (2010) offered guidelines for NSIs adopting either system without demanding that either be adopted.

A non-systematic poll of 12 RDCs (ADSS, 2016) found that RDCs were 50-50 split between rules-based and principles-based OSDC. However, discussions of the merits of the two are largely confined to practitioner meetings or papers; for example, Lowthian and Ritchie (2017) discuss how principles-based operates in an academic research network. The only peer-reviewed paper (Ritchie and Elliot, 2016) directly addressing the topics is from the principles-based camp. Ritchie and Elliot (2016) examine the PB/RBOSDC debate, arguing strongly that the principles-based system is superior; however, they acknowledge that the principles-based model requires a greater institutional commitment, and that the rules-based model is an easier ‘sell’ to the data holders.

Finally, in 2017 the UK Office for National Statistics (ONS) revised the national training for UK-based researchers working with confidential microdata (ONS, 2019). The previous training model, which dominated UK training from 2004 and strongly influenced other countries’ confidentiality training, treated OSDC as a statistical problem. The revised model was the first document to be explicit about the operational justification.

### Models of user segmentation

PBOSDC implicitly acknowledges that research and researchers have multiple skills, interests and demands. As Ritchie (2007) notes, this problem becomes manageable when considering how demand inputs can be segmented. The notion that different types of requests from customers require different approaches to operational delivery is well-established in the discipline of management.

The foundations of this approach can be identified in research on improving operational efficiency. Whilst exploring methods of increasing effectiveness of Just-in-Time (JIT) manufacturing strategy, Pareto analysis was applied to manufactured products to describe the demand pattern of products originally identified as “regular runners, irregular runners, and strangers” (Parnaby, 1988: 486). The categorisation was used to better understand the predictability of the customer request and its impact on availability of organisational resources required to fulfil the order. Parnaby proposed that efficiency gained through JIT success relied on a dependable stream of resources for ‘runners’ and ‘irregular runners’ (later called ‘repeaters’). ‘Strangers’ require increased levels of customised work, making it less amenable to JIT workflow management and therefore less efficient.

While Parnaby does not define these labels, the terms are described in a seminal Business Process Management (BPM) paper by Armistead (1996).

- Runners – demand which is part of the regular routine, predictable resource requirement
- Repeaters – intermittent and uncertain demand, some known resource requirement
- Strangers – much less predictable demand, very limited insight for resource allocation

‘Aliens’ were a later addition (Aitken et al., 2003) describing requests from the customer which are so infrequent or unfamiliar that pre-existing knowledge is generally not applicable. Thus, a state of ‘readiness’ for forecasting resources for such a request cannot be achieved.

Armistead (1996) draws attention to the connection between variety in customer demand and the resource consumed in the production process. In his view, demand variety has multiple dimensions: changes in

volume and differences in requested output. This connection draws heavily on a concept especially relevant here, Ashby's (1956) 'Law of Requisite Variety'. Requisite variety mandates that any system must meet request variety with a similar variety in production capability; or it must attenuate/reject that request to remain viable. Thus, the success or failure of a delivery system is determined by its adequacy in managing its environment of customers and suppliers (Pickering, 2002; Beer, 1984).

The categorisations of demand characteristics act as an aid to the organisation in managing its environment and maintaining viability through the efficient allocation of resource. In this way, efficiency can be seen as the product of how well the delivery process is designed to meet the variety in demand.

Alignment between the design and the context in which it will operate has been shown to lead to optimal performance (Frei, 2006, Sampson & Froehle, 2006). Similarly, research has identified a connection between design and performance, whereby "inadequate service design will cause continuous problems with service delivery" (Gummesson, 1994: 85). Considering the potential applications in the context of the ONS, research by Sousa & Voss (2006) may be highly relevant: in the face of higher request variety, an organisation can employ a design strategy which uses different operational means of delivering similar outputs to customers.

This concept was empirically explored in Ponsignon et al. (2011) where complexity of customer demand was shown to determine the level of customisation provided by the delivery system. This approach provides benefits from efficiency created through standardisation for 'runners', while enabling the organisation to react to complex inputs with customisation for the 'strangers'. The unfamiliar nature of 'aliens' may require innovation in process design in order to accept the related presented variety.

Encountering 'strangers' and 'aliens' forces an organisational choice of whether to accept the input request, or attenuate the variety and reject the request. If accepted and produced, the new output may then be offered to other customers by continued implementation of the newly-created process (Aitken et al., 2003).

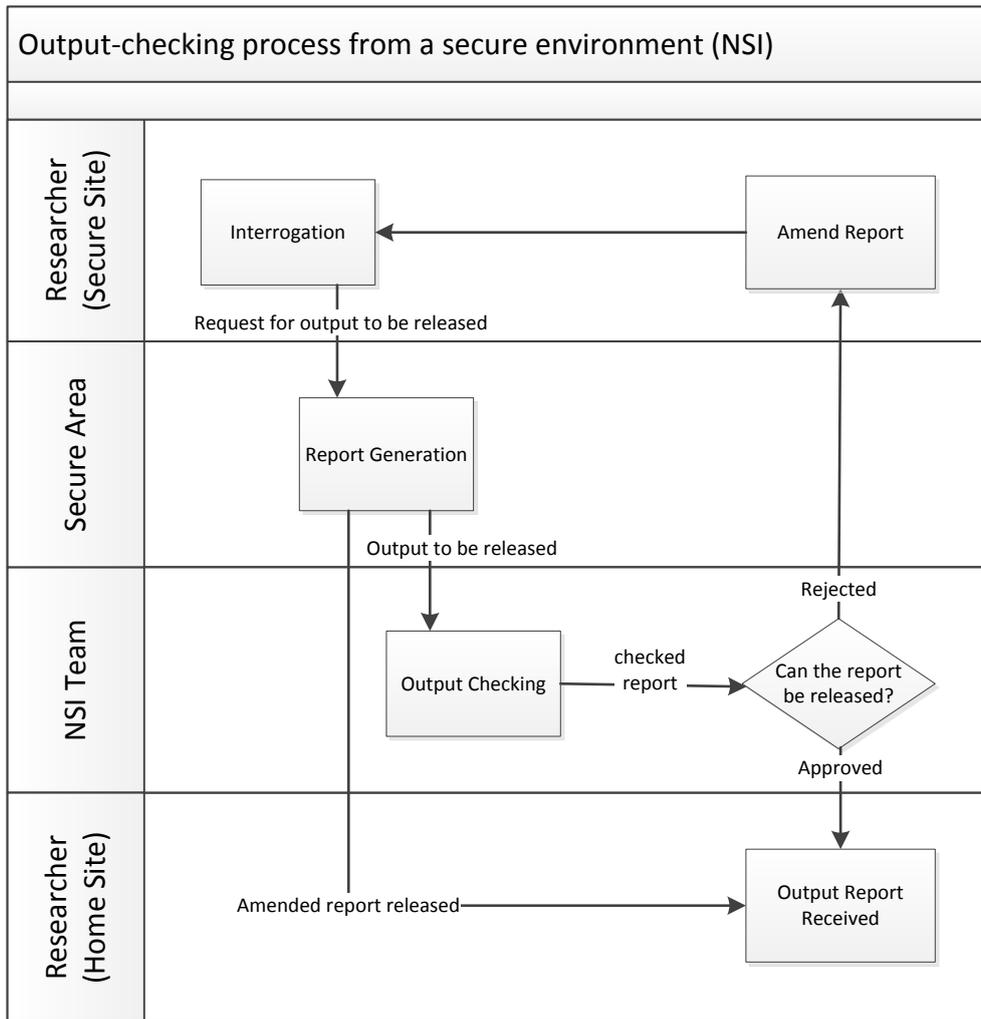
Conversely, the organisation may implement design which requires greater participation by the customer in the creation of the output. Frei (2006) suggests the accommodation of customer-presented complexity through 'low-cost accommodation'. By shifting work away from the organisation and back to the customer, the organisation can derive some benefit from efficiencies in resource allocation. In this case, customers are given access to the delivery system in order to 'self-serve' and create their own outcomes.

Sufficient evidence exists to support the application of the RRSA model to OSDC for the purposes of increasing efficiency in the use of resources through adjustments to the organisational delivery system. Central to this, it is necessary to explore the alignment between the nature of the request from the customer and the process required to fulfil that request.

### 3. Rules-based, principles-based and ad-hoc output-checking

Figure 1 below shows a typical output-checking process from a secure environment managed by a National Statistical Institute (NSI):

Figure 1 Example output-checking process



The researcher works in an environment where he or she cannot directly take away statistical results (note: some facilities allow more ‘trusted’ users to check and release their own outputs). The researcher places the outputs to be released in some predefined location in the secure environment and asks the support team to check and release the output. The support team can extract results from the secure environment. If the support team decides the output is non-disclosive, it sends the results out to the researcher’s (open) home environment.

For expository purposes, we will assume that the researcher has asked for a frequency table to be released, and that the support team operates a simple threshold rule of three; that is, the table must have at least three observations underlying each cell in the table. So, in the example below, Table (a) passes the SDC rule but table (b) does not:

(a) Age versus diabetic status				(b) Gene marker vs diabetic status			
		Age group				Genetic marker	
		18-24	25-29			Yes	No
Men	Diagnosed	11	9	Men	Diagnosed	18	2
	No diagnosis	349	407		No diagnosis	72	684
Women	Diagnosed	12	14	Women	Diagnosed	21	5
	No diagnosis	267	299		No diagnosis	64	502

Note: all data fictional and for illustrative purposes only

Under a rules-based approach (RBOSDC), this is a hard limit; no exceptions are allowed. Under the principles-based approach (PBOSDC), the researcher can argue that the rule is inappropriate in the following circumstances (ONS, 2019):

- if the output is non-disclosive, and
- if the detail in the output is important to the researcher, and
- if this request for an exception is a rare occurrence for the researcher

The first condition is the obvious minimum. The second condition ensures that the output-checker and researcher only spend time negotiating over an output when the result matters to the researcher. This is appealing to researchers as it puts them in charge of deciding when something is ‘important’, rather than the output checker. Thus, Table (b) above could be released if the researcher demonstrated that the small value was non-disclosive and essential for publication. The third condition ensures that researchers do not abuse the system. Note that the terms “important” and “rare” are not specified – this is an area for the researcher and output-checker to negotiate (ONS, 2019). As a result, training the researcher to understand the concept is necessary; effectively, this is Frei’s (2006) model of low-cost accommodation.

Under PBOSDC, the output-checker can also argue that the rule is inappropriate in a specific case because it does not protect confidentiality. For example, in the above case the output-checker may argue that a higher threshold is needed because the data are particularly sensitive and the patients are easily identified. Some organisations (for example, National Records for Scotland) operate a two-tier system with a lower ‘regular’ threshold and a higher threshold for outputs based on Census data.

To give more certainty to the researcher over what will or will not be allowed, PBOSDC systems usually use higher thresholds (10 is common) than RBOSDC. Use of an overly-restrictive rule should not limit research as the researchers always have the opportunity to argue for an exception. In other words, the ‘rule’ is now a rule-of-thumb which can be adjusted up and down as necessary; the rule-of-thumb is designed to be ‘good enough’ in most circumstances; and it can be set much more strictly than in the rules-based case because there is always the option to adjust when important.

This combination of stricter rules-of-thumb and the ability to use discretion in applying those stringent rules is what gives the principles-based approach its superior risk management. Under RBOSDC, a single rule has to do two jobs: protecting confidentiality (by having a higher threshold, for example), and allowing useful, non-disclosive output to be published (which is limited by having a high threshold). Security and efficiency must be traded off. In contrast, under PBOSDC, the rule has one job (protect confidentiality in most cases); efficiency comes through negotiation when it matters.

Rules-based models also fail to provide the imagined guarantees over security. Consider the following tables:

(c) Proportion with no genetic markers				(d) Diabetes diagnosis versus BMI					
		Number	No genetic Marker		Body mass index				
					<18	18-25	25-30	>30	
M	Diag.	20	90%	M	Diag.	0	0	3	17
	No diag.	756	10%		No diag.	110	511	94	41
F	Diag.	26	81%	F	Diag.	3	3	4	16
	No diag.	566	11%		No diag.	46	449	56	15

In Table (c), all cells have at least 5 underlying observations. However, it is clear that an *implicit* table is being generated: the complement to the proportion with the genetic marker is the proportion without it. Table (c) shows that there are 2 males (10% of the 20 in total), diagnosed with diabetes in the dataset who have the genetic marker.

Table (d) shows the problem of *class disclosure*. All males in this dataset diagnosed with diabetes have a BMI greater than 25 i.e. they are overweight or obese. It doesn't matter that there are twenty individuals in this group, well above the threshold; something is now known about *all* males with diabetes diagnosed.

This is the simplest statistical case for PBOSDC over RBOSDC; other examples can be developed. When combined with the higher thresholds used in PBOSDC, it is clear that RBOSDC is the higher-risk option. If this is the case, why do risk-averse organisations use RBOSDC? Two reasons are invariably given.

First, rules are said to be simpler for everyone to use (researchers and output-checkers) and easier to explain to data-holders who want to be reassured when depositing their data. The latter is a valid point: data holders, if aware of SDC at all, are likely to be familiar only with the traditional model of SDC for tabular data in a hostile environment. Simple rules reflecting that knowledge have immediate appeal, even though the sense of security in the familiar is not warranted.

The second, and more common, reason given is that rules-based uses fewer resources: applying simple rules should be easier and require lower-skilled operators than a system which leaves open the possibility of negotiation over any statistical artefact. Principles-based systems cannot be less resource-intensive than rules-based models in the absence of queries, and must require more resources if the checking staff must deal with queries. Moreover, those resources involve output-checkers with statistical skills, which are not necessary for the rules-based system.

There is a third option which is widely implemented. Almost no rules-based organisations operate in the simplistic way described above. All have some informal arrangement allowing researchers to argue that outputs which break the rules can be released in certain circumstances. We will refer to this as 'ad hoc' output SDC (AHOSDC). This method can provide some of the flexibility/efficiency gains of the full principles-based approach without the potential free-for-all. At first glance, this approach seems to offer the best of both worlds.

In practice, it suffers the key problems of both. First, it does not address the risky nature of rules-based by making no allowance for output checkers ignoring rules to block disclosive outputs. More importantly, not formally acknowledging that rules are flexible can create a lack of clarity, causing uncertainty and inefficiency. It also makes resource allocation harder: should output checkers have statistical skills when the formal policy of the organisation says that they do not need them?

Some organisations argue that the simple threshold rule presented above is a straw man: more complicated rules can achieve both the security and the flexibility of PBOSDC. Unfortunately, this is extremely difficult to do in genuine research environments. Ritchie (2007) provides a counter example where a simple, specific, unambiguous, 17-word threshold rule rapidly becomes a woolly 47-word mouthful which requires specialist interpretation, and which is easily challenged by a researcher wanting to make a point.

Fundamentally, the reason why RBOSDC (and AHOSDC) fails to meet the twin targets of efficiency and security is because it is grounded in the SDC literature which sees this as a statistical problem generated by an arbitrary 'user' type of individual. In contrast, PBOSDC sees output checking as a process problem, caused by multiple types of client and client needs. To see why this makes such a difference in implementation, we now turn to the management literature.

#### 4. Output checking as a user segmentation problem

The process perspective on output checking starts from the recognition that different types of researchers, and types of output, produce different demands on the NSI. As noted in the literature review, the concept of 'requisite variety' was established as far back as 1958, and there is a well-established management literature which uses segmentation of customer demand as a way of efficiently allocating resources. Simpler demands are automated as far as possible, leaving specialist resources to be concentrated on the more specialist, high-value cases.

We employ the ‘runners, repeaters strangers and aliens’ (RRSA) model. Using the RRSA terminology, we can divide output requests into four separate types, with some indication of how often these occur

Type (frequency in outputs)	Characteristics	Example	Resource need
Runners 80%-90%	Outputs that could be checked automatically	Small simple tables exceeding rules of thumb; regression coefficients; concentration indexes	Checks for classification and yes/no rules, with an assumption of clearance
Repeaters 10%-20%	Outputs that require human but non-technical review	Multiple linked tables, large or multidimensional tables; graphs; tables where the numbers fall below the threshold	Simple tests applied to provide assurance; assumption of clearance given context
Strangers 1%-2%	Outputs requiring technical review and the development of new guidelines	New statistical outputs with no current guidelines; datasets with very unusual characteristics	Detailed review by technical staff plus development of new guidelines
Aliens n/a	Outputs not normally considered as relevant to this environment	Release of record-level data rather than statistics; release of qualitative data e.g. quotes or video images	Review of appropriateness of environment

The runners are the bread-and-butter of microdata research. They include simple descriptive statistics such as mean of the observations, or frequency counts in categories, which are usually presented with high numbers of observations. The runners also include ‘safe’ (or ‘low review’) statistics, such as regression coefficients, where there is no meaningful disclosure risk (Ritchie, 2019). These outputs could in theory be reviewed automatically, using for example the programs tau-Argus or sdcMicro; in practice they are manually reviewed as this is faster.

The repeaters are the outputs which require the reviewer to make a judgement based on context. For example, a scatter plot of regression residuals might be submitted; the checker would want to evaluate the risk in any outliers. Alternatively, this could be a simple table with counts below the threshold (as in a PBOSDC ‘exception’ request): the checker is then being asked to make a judgment on whether this is non-disclosive, infrequent and important.

These two cover almost all outputs from research centres. Note that in RBOSDC, only the runners exist: an output cannot be cleared unless a known unambiguous rule exists. For PBOSDC, allowing for repeaters is essential: this flexibility to review outputs in context allows much more restrictive (that is, protective) rules to be placed on runners.

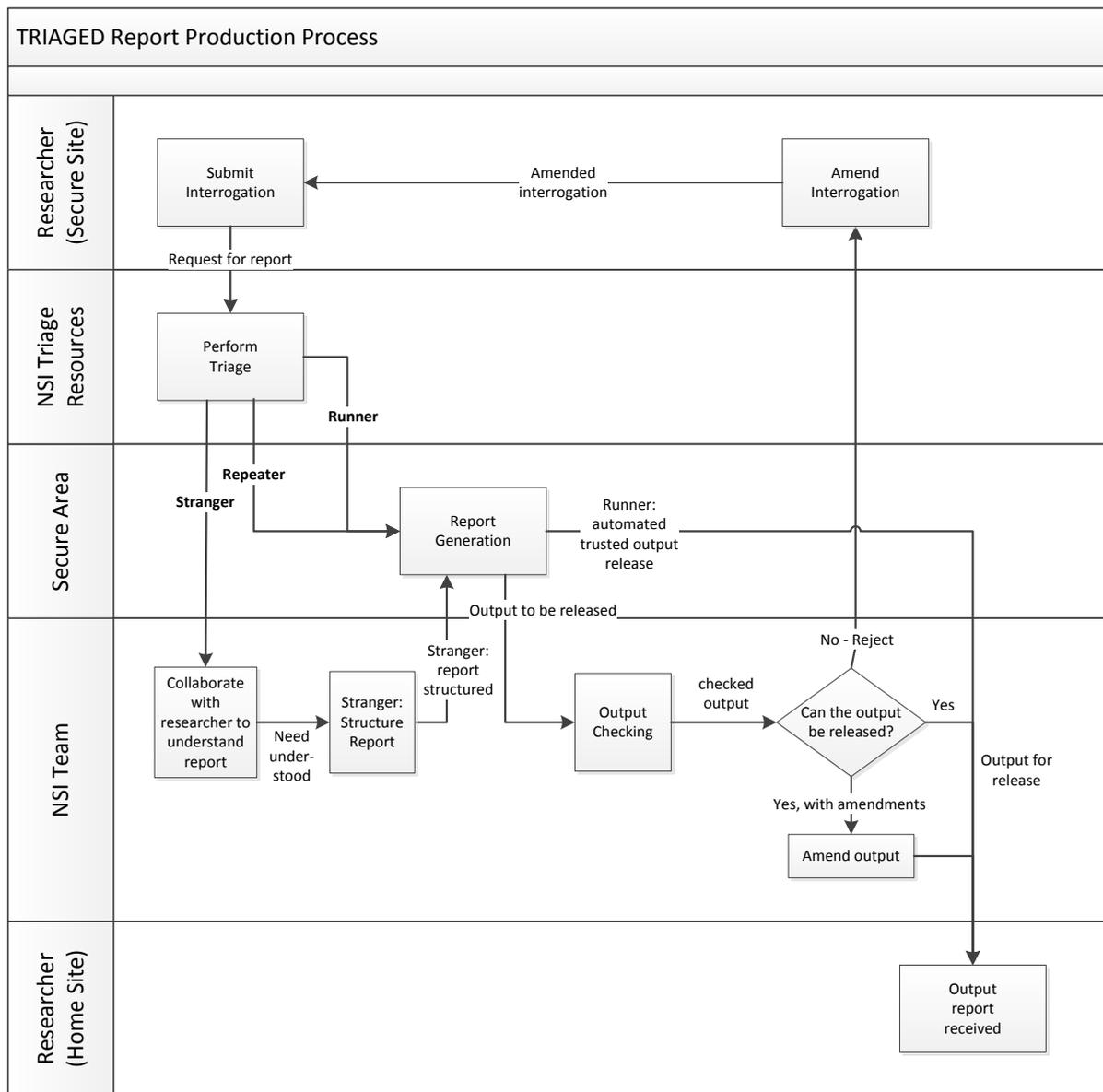
The strangers are where the researcher produces something that the output checking team hasn’t seen before. It could be a novel output (for example, being asked to make a decision on a Herfindahl index for the first time), or familiar outputs presented in an unfamiliar way (in one case, a project which required a very large number of intersecting tables). These require multiple skills: a reasonable degree of statistical knowledge, an ability to judge evidence effectively, and the social skills to hold productive discussions with the researcher.

This highly-skilled resource is expensive; analytical work is the most efficient utilisation of that resource, and so limiting the time spent by that resource on checking outputs is important for the organisation. Ideally, a stranger output only appears once: the role of the reviewer is to decide whether this specific output is to be released, and how future outputs of the same type should be classified. For example, on first encountering a heat map or box-and-whisker plot, the former would be classified as a runner, the latter as a repeater.

Finally, the aliens are those outputs for which the facility was not designed; for example, the release of a linked dataset constructed by a researcher. This does not require statistical knowledge at all, but rather understanding of the purpose of the facility. It may lead to a redesign of the facility (say, an isolated section to allow linking to take place).

An amended model of the process is presented in Figure 2, which reflects the various process flows associated with a triage activity sorting runners, repeaters, and strangers.

Figure 2 Output checking when viewed as a multi-stage triaged process



As well as managing resources, this structure also allows the NSI to make the researcher an active part of the clearance process; an approach which utilises Frei’s (2006) ‘low cost accommodation’. ONS’s (2019) PBOSDC training for researchers emphasises three points:

- Runners are done quicker than repeaters
- If your planned output is a stranger, get the review team involved as soon as possible, not when you want the output; then you don’t waste time producing unreleasable outputs
- Provide all the information for the reviewer to put the output into one of the classes

The aim is to make the researcher see that his or her behaviour has a direct impact on clearance times, and to show the researcher how to improve them. By making the researcher an active part of the clearance chain, carrying out the preparatory work, the output-checker finds his or her workload reduced. The

awareness of researchers that they can directly affect response times builds a feeling of control and hence engagement in the process.

The RRSA model also provides a clear structure for staff resources. Consider the skills needed by an output-checker for the different types of output, and that person’s discretion to ignore the rules of thumb:

Type	Skills	Rules of thumb	Standardisation
Runners	Ability to recognise types of output and follow rules	Follow	Highly standardised
Repeaters	Good understanding of data and practical (not theoretical) understanding of disclosure risk; statistically competent but not expert	Follow with interpretation	Mixed standardisation and discretion
Strangers	Statistical/data skills to understand new types of problems and take decisions	Develop new ones	High use of discretion
Aliens	Strategic perspective on operations	Out of scope	n/a

This creates a hierarchy of technical skills allowing different staff to be allocated to different roles. It also simplifies skills acquisition and staff training, by providing a clear path to personal development based upon experience and knowledge of the data.

This differentiation of skills is the second reason why apparently lower-cost models of output clearance fail to achieve the operational gains of PBOSDC. For RBOSDC, only runners and strangers should exist: there are fixed rules, which may be added to as new statistical products occur. This implies that the bulk of the work can be carried out by checkers with minimal training in statistics or data.

However, the systems run by most NSIs do not follow a hard-rules model but are ad hoc (AHOSDC); that is, notionally the rules are hard but in practice researchers ask for, and get, some flexibility. The flexibility may depend on the data, the statistic, or sometimes whether the researcher is ‘trusted’ or not. The flexibility is important: without it, the NSI is likely to lose the goodwill of the researcher. The difficulty is that, because the flexibility is not officially sanctioned, it is less clear whether a clearance is going to be simple or complex, allowed or blocked. Clearance times become less certain; and because any clearance might be an exception, all clearance staff need to have the ability to handle exceptions. The efficiency gains from having clearly delineated production processes have been lost.

## 5. Cost-effectiveness and resource use

PBOSDC was devised and first implemented at the UK Office for National Statistics (ONS). While the statistical benefits became evident over time, the initial appeal was as a way of keeping overall costs down via low-cost accommodation. In this, it appeared to be successful: at its initial peak in 2008-2010, ONS’ secure research facility was reportedly releasing more outputs from more researchers at lower staff cost than comparable European facilities<sup>1</sup>. The ONS model also scaled easily: in 2010, some 2000 release requests (a request could be anything from a single regression to tens of linked tables or graphs) were easily handled by one statistically competent full-time equivalent (FTE)<sup>2</sup>.

However, there are extremely resource-efficient rules-based systems. Statistics Norway runs both a full-service RDC, and a remote job model (microdata.no) developed in collaboration with the Norwegian Center for Research Data (NSD)<sup>3</sup>. Microdata.no takes the rules-based model to its logical conclusion: all decisions are taken by computer, and the system does not allow outputs which do not have a clearance rule attached. This

<sup>1</sup> This information was gained in conversations for Eurostat expert group 2009-10, and from presentations by Scandinavian and North American RDC operators.

<sup>2</sup> The team actually had five output-checkers, each of whom spent one day a week checking outputs. In practice they reported spending 2-3 hours on their allotted day, implying rather less than one FTE.

<sup>3</sup> <https://microdata.no/>

is highly resource-efficient and allows users to see release decisions in real time. As a result, despite the very strict rules to manage disclosure risk (for example, the initial threshold for tables is set at minimum of 1,000 observations), user responses have been very positive and the model has attracted substantial interest from NSIs and other organisations.

It is also feasible to run very cost-effective ad hoc systems. In social science, the pre-eminent example is LISSY<sup>4</sup>, which has been running for almost two decades and allows users to submit code to run analyses of the Luxembourg Income/Wealth Studies. Like microdata.no, simple strict rules are applied automatically by the server (including allowed commands), but the computer's triaging allows for the option "set for review" (that is, send to a human for checking). As a result the rules are less stringent than microdata.no. As in Norway, users get immediate feedback on whether code will be allowed to run or not, and users are encouraged to recode rather than waiting for review. Despite the small staff, LISSY handled 73,000 jobs in 2018, and has shown continual growth in both user numbers and data requests, indicating a high level of user satisfaction.

Thus, there are examples of efficient principles-based, rules-based or ad hoc OSDC systems. Methodological problems limit the chance of comparative evaluation, but it is clear that there is no solid evidence to support the argument that PBOSDC is more expensive than other solutions and not scalable. On the contrary, seeing the problem from a management perspective makes clear that we should expect PBOSDC to be more efficient than ad hoc solutions, more flexible than rules-based solutions, and easily scalable as long as the investment in training researchers is made.

Finally, it is worth applying the RRSA model to the other examples above to show how this perspective helps us understand their efficiencies too. The Norwegian model only has runners and aliens; the latter are used to identify new rules to widen the class of runners. With only runners, automatic clearance of all outputs is the logical and cost-effective conclusion. In the case of LISSY, repeaters are allowed but strongly discouraged via immediate feedback; and just as for PBOSDC, this feedback is designed to encourage users to change their behaviour.

## 6. Conclusion

Secure research access to the most sensitive microdata has been one of the great success stories for NSIs this century. It came from realising that simply reducing data detail was a dead-end; instead, novel ways of working with researchers opened a range of options. For all of these new ways of working, output-checking is a key part of the operational system.

Perceptions of output-checking have been dominated by the statistical literature, which is designed to address the safe production of statistical aggregates. Statistical aggregates are well suited to a rules-based system, but research outputs are not. Hence OSDC was born as a field, with PBOSDC its standard-bearer. But to those brought up on the traditional statistics, PBOSDC seemed an operational nightmare: how can an explicitly 'flexible' (i.e. uncertain) world, requiring greater statistical understanding and more training for everyone, be both safe and scalable?

When viewed from an operations management perspective, the answer is that the efficiency gains come precisely from the elements that worry traditionalists. Using crude rules which have a large margin of error but with flexibility at the margin for high-value outputs means that the 80%-90% of 'runners' can be handled quickly by automatic process or staff with minimal training. Allowing researchers to choose when their runners become 'repeaters' saves the output-checker carrying out this function. This is low-cost accommodation: the NSI has effectively turned the customer into part of the workforce. More importantly, it gives the researcher some control over the process, and so builds engagement. There are upfront training costs, but these should be seen as investment expenditure.

---

<sup>4</sup> <https://www.lisdatacenter.org/data-access/lissy/>

From the management literature, there are no surprises that a one-size-fits-all model allocates resources less efficiently than a segmented-markets model; nor that the latter is better at exploiting customer self-service. This 'requisite variety' has been a core of management thinking for over half a century. What is perhaps less obvious is that this also produces better statistical outcomes: PBOSDC is inherently lower-risk than RBOSDC (basic rules are stricter; resources are concentrated on checking high-risk outputs). It also reduces dissatisfaction amongst users, a known risk factor for restricted-access systems.

This illustrates a wider issue. The traditional focus on statistical measures of risk, without considering the implications of operational choices, has been strongly criticised (e.g. Hafner et al., 2015) as risky and inefficient; Ritchie and Smith (2018) also suggests that this is doomed to failure in a big data/machine learning world. In contrast, operations research has much to say about effective risk management, particularly in relation to digital services (such as the 'data supply chain' model of Spanaki et al., 2017). A change in emphasis from statistical to operational models of risk drawing on the extensive management literature (as in ONS' 2019 course for output checkers which uses the SSRA framing), should help NSIs to improve delivery on the joint objectives of security, efficiency, and customer service.

## References

- ADSS (2016) Data Access Strategy: final report. Australian Department of Social Services, June.
- Aitken, J., Childerhouse, P. and Towill, D., 2003. The impact of product life cycle on supply chain strategy. *International Journal of Production Economics*, 85(2), pp.127-140.
- Armistead, C. (1996). Principles of business process management. *Managing Service Quality: An International Journal*, 6(6), 48-52.
- Ashby, W. R. (1956), *An Introduction to Cybernetics*, Chapman & Hall Ltd., London.
- Beer, S. (1984). The viable system model: Its provenance, development, methodology and pathology. *Journal of the operational research society*, 35(1), 7-25.
- Bond S., Brandt M., de Wolf P-P (2015) Guidelines for Output Checking. Eurostat.  
[https://ec.europa.eu/eurostat/cros/system/files/dwb\\_standalone-document\\_output-checking-guidelines.pdf](https://ec.europa.eu/eurostat/cros/system/files/dwb_standalone-document_output-checking-guidelines.pdf)
- Brandt M., Franconi L., Guerke C., Hundepool A., Lucarelli M., Mol J., Ritchie F., Seri G. and Welpton R. (2010), Guidelines for the checking of output based on microdata research, Final Report of ESSnet Sub-group on Output SDC [http://neon.vb.cbs.nl/casc/ESSnet/guidelines\\_on\\_outputchecking.pdf](http://neon.vb.cbs.nl/casc/ESSnet/guidelines_on_outputchecking.pdf)
- Chase, R. (1981) "The Customer Contact Approach to Services..." *Operations Research*, 29, 4.
- Hafner H-P., Lenz R., Ritchie F., and Welpton R. (2015) Evidence-based, context-sensitive, user-centred, risk-managed SDC planning: designing data access solutions for scientific use. In: *Worksession on Statistical Data Confidentiality 2015*, Eurostat.
- Hundepool, A., Domingo-Ferrer, J., Franconi, L., Giessing, S., Lenz, R., Longhurst, J., Schulte Nord-holt, E., Seri, G. and De Wolf, P-P. (2010). Handbook on Statistical Disclosure Control. ESSNet SDC.  
[http://neon.vb.cbs.nl/casc/.SDC\\_Handbook.pdf](http://neon.vb.cbs.nl/casc/.SDC_Handbook.pdf)
- Corscadden, L., Enright J., Khoo J., Krsnich F., McDonald S., and Zeng I. (2006) Disclosure assessment of analytical outputs. Mimeo, Statistics New Zealand, Wellington.
- Eurostat (2016) Self-study material for the users of Eurostat microdata sets.  
<http://ec.europa.eu/eurostat/web/microdata/overview/self-study-material-for-microdata-users>
- Frei, F. X. (2006), "Breaking the Trade-Off Between Efficiency and Service", *Harvard Business Review*, 84, 11, pp. 92-101.

- Lowthian P. and Ritchie F. (2017) *Ensuring the confidentiality of statistical outputs from the ADRN*. Technical report no3. Administrative Data Research Network
- O'Keefe C., Westcott M., Ickowicz A., O'Sullivan M. and Churches T. (2015) *Guidelines for Confidentiality Protection in Public Health Research Results*. CSIRO.
- ONS (2019) *Safe Researcher Training [2017 onwards]*. Office for National Statistics, Research Support Service. Last viewed June 2019.
- Parnaby, J. (1988). A systems approach to the implementation of JIT methodologies in Lucas Industries. *The International Journal Of Production Research*, 26(3), 483-492.
- Pickering, A. (2002). Cybernetics and the mangle: Ashby, Beer and Pask. *Social studies of science*, 32(3), 413-437.
- Ponsignon, F., Smart, P. A., & Maull, R. S. (2011). Service delivery system design: characteristics and contingencies. *International Journal of Operations & Production Management*, 31(3), 324-349.
- Reznek, A. (2004) *Disclosure risks in cross-section regression models*, mimeo, Center for Economic Studies, US Bureau of the Census, Washington
- Reznek A. and Riggs T. (2005) *Disclosure risks in releasing output based on regression residuals*. ASA 2004 Proceedings of the Section on Government Statistics and Section on Social Statistics pp1397-1404
- Ritchie F. (2006) *Disclosure control of analytical outputs*. Mimeo: Office for National Statistics. Edited and reprinted as WISERD Data and Methods Working Paper no. 5 (2011).
- Ritchie F. (2007) *Statistical disclosure control in a research environment*. Mimeo, Office for National Statistics. Edited and reprinted as WISERD Data and Methods Working Paper no. 6 (2011).
- Ritchie F. (2008) *Disclosure detection in research environments in practice*. In: *Work session on statistical data confidentiality 2007*, Eurostat; pp399-406
- Ritchie F. (2014) *Operationalising safe statistics: the case of linear regression*. Working papers in Economics no. 1410, University of the West of England, Bristol. September
- Ritchie F. and Elliot M. (2015). Principles- versus rules-based output statistical disclosure control in remote access environments. *IASSIST Quarterly* 39:5-13
- Ritchie F. and Smith J. (2018) *Confidentiality and linked data*. In *Statistics Authority (2018)*
- Sampson, S. E. and Froehle, C. M. (2006), "Foundations and Implications of a Proposed Unified Services Theory", *Production and Operations Management*, 15, 2, pp. 329-343.
- Spanaki, K., Gürgüç, Z., Adams, R., & Mulligan, C. (2018). Data supply chain (DSC): research synthesis and future directions. *International Journal of Production Research*, 56(13), 4447-4466.
- Sousa, R. and Voss, C. A. (2006), "Service Quality in Multichannel Services Employing Virtual Channels", *Journal of Service Research*, 8, 4, pp. 356-371.
- Statistics Authority (2018) *National Statistician's Quality Review on Privacy and Confidentiality*. Ed. G. Roarson. UK Statistics Authority, December.
- Stats NZ (2015). *Microdata Output Guide (Third edition)*. Statistics New Zealand. Available from [www.stats.govt.nz](http://www.stats.govt.nz).
- Willenborg L. de Waal T. *Statistical Disclosure Control in Practice: v. 111*. Springer: New York: Lecture Notes in Statistics; 2013

## **Recent UWE Economics Papers**

See <https://www1.uwe.ac.uk/bl/research/bcef/publications.aspx> for a full list.

### **2019**

- 1904 **Runners, Repeaters, Strangers and Aliens: Operationalising efficient output disclosure control**  
Kyle Alves  
Felix Ritchie
- 1903 **Artificial Intelligence and the UK Labour Market: Questions, Methods and a Call for a Systematic Approach to Information Gathering.**  
Tim Hinks
- 1902 **Robots and Life Satisfaction**  
Tim Hinks
- 1901 **Education and the Geography of Brexit**  
Robert Calvert Jump and Jo Michell

### **2018**

- 1808 **Pricing behaviour and the role of trade openness in the transmission of monetary shocks**  
Laura Povoledo
- 1807 **Learning, Heterogeneity, and Complexity in the New Keynesian Model**  
Robert Calvert Jump, Cars Hommes, and Paul Levine
- 1806 **DSGE Models and the Lucas Critique. A Historical Appraisal**  
Francesco Sergi
- 1805 **A new approach to estimating interregional output multipliers using input-output data for South Korean regions**  
Malte Jahn, Anthony T. Flegg and Timo Tohmö
- 1804 **Urban food security in the context of inequality and dietary change: a study of school children in Accra**  
Sara Stevano, Deborah Johnston and Emmanuel Codjoe
- 1803 **The use of differential weighting and discounting in degree algorithms and their impact on classification inflation and equity: A further analysis**  
David O. Allen
- 1802 **Unambiguous inference in sign-restricted VAR models**  
Robert Calvert Jump
- 1801 **Degree algorithms, grade inflation and equity: the UK higher education sector**

David O. Allen

## **2017**

- 1706 **Internal rationality, heterogeneity and complexity in the new Keynesian model**  
Cars Hommes, Robert Calvert Jump and Paul Levine
- 1705 **The regionalization of national input–output tables: a study of South Korean regions**  
Anthony T. Flegg and Timo Tohmo
- 1704 **The impact of quantitative easing on aggregate mutual fund flows in the UK**  
Iris Biefang-Frisancho Mariscal
- 1703 **Where are the female CFOs?**  
Gail Webber, Don J Webber, Dominic Page and Tim Hinks
- 1702 **Mental health and employment transitions: a slippery slope**  
Don J Webber, Dominic Page and Michail Veliziotis
- 1701 **SMEs access to formal finance in post-communist economies: do institutional structure and political connectedness matter?**  
Kobil Ruziev and Don J Webber

## **2016**

- 1611 **Curriculum reform in UK economics: a critique**  
Andrew Mearman, Sebastian Berger and Danielle Guizzo
- 1610 **Can indeterminacy and self-fulfilling expectations help explain international business cycles?**  
Stephen McKnight and Laura Povoledo
- 1609 **Pricing behaviour and the role of trade openness in the transmission of monetary shocks**  
Laura Povoledo
- 1608 **Measuring compliance with minimum wages**  
Felix Ritchie, Michail Veliziotis, Hilary Drew and Damian Whittard
- 1607 **Can a change in attitudes improve effective access to administrative data for research?**  
Felix Ritchie
- 1606 **Application of ethical concerns for the natural environment into business design:**

**a novel business model framework**

Peter Bradley, Glenn Parry and Nicholas O'Regan

- 1605 **Refining the application of the FLQ Formula for estimating regional input coefficients: an empirical study for South Korean regions**

Anthony T. Flegg and Timo Tohmo

- 1604 **Higher education in Uzbekistan: reforms and the changing landscape since independence**

Kobil Ruziev and Davron Rustamov

- 1603 **Circular economy**

Peter Bradley

- 1602 **Do shadow banks create money? 'Financialisation' and the monetary circuit**

Jo Michell

- 1601 **Five Safes: designing data access for research**

Tanvi Desai, Felix Ritchie and Richard Welpton

**2015**

- 1509 **Debt cycles, instability and fiscal rules: a Godley-Minsky model**

Yannis Dafermos

1508 **Evaluating the FLQ and AFLQ formulae for estimating regional input coefficients: empirical evidence for the province of Córdoba, Argentina** Anthony T. Flegg, Leonardo J. Mastronardi and Carlos A. Romero

1507 **Effects of preferential trade agreements in the presence of zero trade flows: the cases of China and India**  
Rahul Sen, Sadhana Srivastava and Don J Webber