



BANK OF ENGLAND
PRUDENTIAL REGULATION
AUTHORITY

Solvency II: internal model approval process data review findings

February 2016

Prudential Regulation Authority
20 Moorgate
London EC2R 6DA

Prudential Regulation Authority, registered office: 8 Lothbury, London EC2R 7HH.
Registered in England and Wales No: 07854923

© Prudential Regulation Authority 2016

Contents

1	Introduction	4
2	Data review: approach	4
3	Summary	6
4	Findings	6
5	Maintenance regime	16
Appendices		17

1 Introduction

1.1 In September 2012, the Financial Services Authority (FSA) published its interim data review findings (the '2012 report') based on reviews of firms in the internal model approval process (IMAP) to assess the quality of data that underpins the validity and integrity of the internal model¹. The report provided feedback to firms based on findings that included an outline of the FSA approach to the data review, a summary of the results, and suggestions for what firms might want to consider when preparing their application to use an internal model.

1.2 The Prudential Regulation Authority (PRA) has now completed IMAP data reviews on over 50 firms in preparation for the introduction of Solvency II² and believes this is an appropriate time to issue a final report. This paper completes and supersedes the 2012 interim report. It includes the findings of the original 2012 report where still appropriate, and the findings from the completed data review work.

1.3 The report is intended to inform firms of the findings, to provide examples of the range of practices that exist in the industry, and to provide commentary on how firms can communicate better with the PRA. This document does not create expectations of firms and nor does it seek to interpret directly applicable regulations, the compliance with which is a matter for firms.

1.4 The report is aimed primarily at firms that have received model approval and who are developing plans for future improvement or maintenance of the internal model: but it will also help firms in the pre-approval process. Firms may want to consider areas of good practice, particularly on how to communicate the details of internal model data flow to the PRA.

1.5 As noted in Sam Woods' letter of 15 January 2016³, it is the PRA's aim is to shed as much light as possible on its approach to taking model approval decisions, and to provide insight into the decision-making process.

2 Data review: approach

2.1 The purpose of the data review was for the PRA to gain assurance that firms have internal processes and procedures to ensure the appropriateness, completeness and accuracy of data used in technical provisions⁴, and in the internal model⁵, and that documentation on data processing complies with the relevant sections of Delegated Regulation 2015/35 ('Delegated Regulation') Article 244

2.2 The review was started by the FSA in September 2011 and completed by the PRA in June 2015. It covered all firms in IMAP across both the life and general insurance sectors. The review process consisted of three 'Acts', as follows.

- **Act I** – an initial visit to review the firm's preparedness and to finalise the scope and timing of the external review of data management (the 'FSA/PRA data audit');

1 See FSA 'Solvency II: internal model approval process data review findings'; September 2012; www.fsa.gov.uk/static/pubs/international/sii-imap-data-review-09-12.pdf.

2 Directive 2009/138/EC.

3 www.bankofengland.co.uk/pradocuments/solvency2/edletter15jan2016.pdf.

4 Directive 2009/138/EC Article 82, see also Commission Delegated Regulation 2015/35 hereafter the 'Delegated Regulation' Articles 19-20.

5 See PRA Rules SCR - Internal Model 11.4 transposing Directive 2009/138/EC Article 121(3) and Delegated Regulation Article 231.

- **Act II** – the firm performed the data audit and the PRA conducted a desk-based analysis of the results and agreed with the firm the timescales for the remediation of any issues identified. The PRA discussed with the firm’s audit team the scope, challenges and significant issues and agreed on any follow up necessary to obtain assurance over outstanding management actions; and
- **Act III** - (not started at the time of the 2012 report), an in-depth review for selected firms of any specific areas of concern highlighted in the previous two Acts, including any significant exclusions from Act II. The objective was to test the original Act II audit for:
 - Scope: was the Act II audit sufficiently comprehensive to cover all areas of data where an error could have a material impact on the solvency calculation and reserving process?
 - Design: was the audit sufficiently well designed to identify the appropriate controls to be tested?
 - Execution: was the audit testing robust and thorough, without placing excessive reliance on self-assessment or assurance from business or management?

2.3 The objective of Act III was not in itself to test the quality of the firm’s controls, nor data quality at the firm, but rather to satisfy the PRA that the Act II audit gave reasonable assurance that the firm’s data quality complies with the regulatory standard. The Act III process was as follows.

- Select data flows to sample, based on a discussion with PRA supervisors and actuaries, and with the firm. The criteria included whether the data flow was of material impact, ie such that data errors would lead to material misstatement of capital resources or capital requirement, and whether the flow had points of complex transformations or joins¹.
- Send an information request to the firm, enumerating the risk factors or business lines, and asking for a high-level description of the data flow for each. A flow diagram was helpful, together with a clear and coherent description. The PRA provided sample diagrams as a guideline, if requested.
- When the firm had supplied an appropriate representation of the flow, form a view on the expected control points, based on the likely impact and risk of the transformations the data is subject to.
- Request the firm to assess the Act II audit against the expected control points as determined above.
- Based on this, take a view whether the Act II audit was of sufficient scope, design and execution to give reasonable assurance that the firm’s data quality complied with the regulatory standard.

2.4 The results from Act III were mixed, and prompted remedial work in some cases. Reference numbers in brackets refer to the original numbering of the 2012 report. For example, (2012.4.13) means section 4.13 of the 2012 report.

¹ A ‘join’ is where two or more data sets are combined together using a common attribute, known as a ‘join key’. For example, information about a stock contract is ‘joined’ to price information using a stock ticker. The ticker is the join key. Keys may include bond identifiers, credit ratings, dates, construction codes etc.

3 Summary

3.1 The September 2012 report identified ten findings, which mapped to the five sections of the FSA data audit: i) approach to managing data; ii) implementation of the data policy; iii) understanding of the data used; iv) controls over data quality; and v) IT environment, technology and tools.

3.2 Improvement since 2012 had been steady. Most firms had embedded data governance into business as usual. While complex IT implementations to support data governance were an issue in 2012, most firms had resolved this, sometimes by decommissioning or modifying the original plan. Impact and risk assessment methodology was standard, and firms understood the need for defining thresholds of materiality in order to apply the principle of proportionality that runs through the Solvency II Directive.

3.3 There was still a disproportionate focus on exposure data (policy data, asset data, and referential data) at the expense of valuation data (prices, rates, ratings, longevity, and claims development) and risk data (historical records of valuation data, assumption stresses etc). Many firms were still interpreting 'data directory' to mean 'data dictionary', and were defining data at a level of granularity that would be difficult for end-users to engage with. Firms were still struggling with appropriate ways of classifying data. Firms sometimes found it difficult to assign data ownership as part of their data governance operating model, due to the movement and transformation of data from one system to another.

3.4 The Act III process demonstrated that firms were not able to represent data flow in a clear or precise way, and were sometimes unable to identify the location of appropriate control points for audits to cover. This meant that some firms did not audit appropriate controls over operations where the risk and the impact of error was high. As a result, some firms had to re-perform audits, and other firms had to revisit elements of the audits. In a handful of cases, significant issues arose from the fresh audit work.

3.5 The findings from Act III are discussed in more detail below in sub-risk 3 ('understanding of the data used') and sub-risk 4 ('IT environment, technology and tools'). Suggestions of how firms can better represent data flows are included in the appendices.

4 Findings

4.1 The findings follow the order of the original 2012 report. Findings that are no longer relevant have been omitted.

Sub-risk 1: approach to managing data

Interpreting and applying a specification for the collection, processing and application of data across the firm.

4.2 The Delegated Regulation Article 244(h) does not refer to a 'data policy' but to a 'specification for the collection, processing and application of data'. It refers to Article 231(3)(e), which requires that the data are collected, processed and applied in a transparent and structured manner, based on a specification of the: definition and assessment of the quality of data, including specific qualitative and quantitative standards for different data sets; use and setting of assumptions made in the collection, processing and application of data; and process for carrying out data updates, including the frequency of regular updates and the circumstances that trigger additional updates.

4.3 Ensuring a consistent interpretation and application of the specification across the firm remained a challenge, however.

Sub-risk 2: implementation of data policy

Measuring, analysing and monitoring data quality in business as usual.

Finding 1: data governance

4.4 2012.4.2 noted that '[S]ome firms were unable to articulate what 'accurate', 'complete' or 'appropriate' meant in practice and were therefore unable to assess data quality effectively'. Some firms were still having difficulty with 'appropriate'. A useful rule of thumb is to check for fitness of purpose every time one set of data is joined to another (for example, ratings data by issuer is joined to an exposure in that rating class). For example, should spread data for corporate bonds with a credit rating of BBB be used to assess the risk of a portfolio predominantly consisting of assets rated BBB minus? It is sometimes possible to quantify the fitness of purpose, for example by sourcing a sample of more appropriate valuation or risk data, and comparing the result to that produced by the data currently in use.

4.5 In the 2012 report most firms had a risk committee or a bespoke 'data steering' committee as a data governance body '[H]owever, it was often difficult to determine a level of summarisation of reporting appropriate and useful to such a committee' (2012.4.4). This remained a difficulty. Some firms still had no reporting of data quality issues and remediation plan to their executive committee, and senior management did not always have adequate oversight of data quality. It is helpful if data governance committee meetings include upstream and downstream data stakeholders. Auditors may want to review minutes to check the involvement of all appropriate business areas in remediation decision making.

4.6 In 2012.4.7, data quality reporting on exposure data (policy terms and conditions, and raw asset data) was generally better developed than for valuation data (prices, rates, longevity data etc) and for risk data (spread movements, price history, improvement distributions or stresses etc). This continued to be a challenge for many. One firm had to re-audit its credit data flow after it was discovered that ratings data and risk factor groupings had been omitted from the original audit. There was still considerable reliance on expert judgement for validating risk data, and there was still a perception that actuarial systems and processes were 'too difficult' for data governance. In reality, the quality of credit spread data, ratings data and other information used for valuation and risk formed the most significant part of the internal model. One firm submitted a complex data diagram for its mortgage portfolio, but did not explain how it joined regional index data to property postcodes in order to estimate the current property value from the most recent valuation.

4.7 In catastrophe modelling, valuation data such as vulnerability curves and event simulations are proprietary, and it is not possible to make a direct assessment of data quality. However, it is often possible to use other data (eg historic claims) to make a reasonableness assessment of expected future claims.

4.8 Most firms were using controls assessment to fulfil this need. Some were using data extract and analysis techniques but this was not the norm. 'Process' includes all business as usual systems, not just strategic IT implementations, to give assurance that the existing system of data governance is adequate. Firms sometimes struggled for years to implement or embed IT systems where a simple analysis of extracted files may have provided insights into data quality.

Finding 2: data ownership

4.9 The 2012 report noted that due to the movement and transformation of data from one system to another, most firms found it difficult to assign data ownership as part of their data governance operating model (2012.4.9). Data are typically produced from multiple source systems upstream and used by many users downstream. Some firms did not have a consistent process to communicate upstream system and process changes and its impact to the downstream users and vice versa.

4.10 There continued to be difficulties with applying data management processes consistently, particularly across large firms, and with ensuring that data producers (upstream) and data consumers (downstream) were adopting consistent processes in managing data from source to the internal model. At one firm, the capital modelling team relied on data supplied by claims and premium systems upstream, yet there were almost no data quality checks performed on upstream systems. At another firm, the actuarial and capital modelling function used data from claims systems to produce reserving and capital figures, even though the claims data upstream was incorrect, and relied on poorly controlled manual processes.

4.11 There was a wide range of data governance models implemented by different firms, which can mostly be classified under the 'three lines of defence' model as follows:

- line 1 is typically the end-users of the data;
- line 2 is typically an independent validation function that is close to the data, ideally has some expert knowledge of the data and how it is used, and periodically review it for completeness, accuracy and appropriateness; and
- line 3 is the audit function. Ideally, auditors will have specialist knowledge of how data is operated on for each risk factor, and where appropriate controls need to be. In practice, auditors are not specialists, and sometimes there is great reliance on self-assessment by line 1, or on control descriptions specified by line 1.

4.12 Different models place different reliance on the lines. For example, some firms combine a strong line 1 with a strong line 3, without much reliance on line 2. Other firms rely on a strong line 2, and the PRA has accepted data audits by line 2, on the principle that '[T]he review should be performed by a suitably qualified person who is independent of model design, build, and operation'.¹

4.13 It is important that the data governance function has a broad understanding of the data flow.

4.14 Some firms had informal, undocumented processes, and some were unable to quantify when errors or control failures were to be escalated for timely resolution. In most cases, understanding the operation and its controls does not involve specialist or expert knowledge. When specialist knowledge is required, some firms have found that involving the capital modelling function – the main users of data – is effective. Where the function is not involved, there is a risk of using incomplete, inaccurate or inappropriate data.

4.15 Firms might want to consider whether:

¹ See 'Data Audit: scope and content', FSA, July 2011; www.fsa.gov.uk/static/pubs/international/external_review_scoping_tool.pdf.

- data quality reports have the results of testing from source to the internal model and cover accuracy, completeness and appropriateness checks;
- the roll-out of guidance to upstream and downstream areas of the business has been reviewed;
- data quality reports have been assessed to understand whether they include the results of testing and remediation from across the data flow;
- expert judgements applied to data are appropriately controlled, for example by-
 - specifying any data or assumptions used in making the judgement;
 - defining and logging the judgements; and
 - validating that the expert judgements are applied and approved by an appropriately experienced individual; and
- there is a documented framework specifying the data to which expert judgement can be applied, the effect on the data, the evidence that is needed to support the expert judgement, and a list of appropriate approvers.

4.16 2012.4.11 stated that ‘a firm’s change management process should be able to rely on the data directory to identify provenance and responsibilities across the data flow and enable effective communication between upstream and downstream users’. It is unhelpful if data directories are so granular that it is difficult for end-users to understand the data flow and communicate their requirements effectively.

Finding 3: groups

4.17 The 2012 report noted that complex data transformations between local entity and group require a consistent interpretation and application of group templates, assumptions and standards (2012.4.14). This remained a challenge, and over-reliance on ‘delegated self-assessment’ or ‘self-certification’ was still an issue (2012.4.13). As noted in 2012.4.15 firms should apply consistency in the standards adopted and the metrics used for monitoring and escalation, even though the individual processes at entity level may be different. While there is a need for proportionality, risks with a similar impact and probability, arising in different entities, should not affect group decision making in materially different ways. Where there is a self-certification based governance mechanism (2012.4.16), there should be a strong process for challenging and auditing the self-assessments.

Sub-risk 3: understanding of the data used

Understanding of the overall data flow by data governance. It is sometimes claimed that a data transformation is for ‘actuaries’ or ‘experts’. However, most operations on data fall into a limited number of types (including input, aggregation, joining and calibration) and it is useful to understand which operations have significant impact and risk, so that appropriate controls can be identified, monitored and reviewed.

Finding 4: impact and risk assessment

4.18 Most firms conducted an impact and risk assessment of the data used in the internal model and associated data processes (2012.4.17), and the application of reasonableness and sensitivity testing was common (2012.4.18). Reasonableness tests included using size of

written premium, technical provisions or regulatory capital as a proxy for 'impact'. Sensitivity testing usually involves changing key data items on the capital path to assess the impact on the capital requirement or technical provisions.

4.19 Applying the concept of materiality consistently (2012.4.19) remained a challenge. It was often difficult to identify:

- risk modules or products that could become material in the future due to changes in the current business model or risk profile;
- instances where a data error could affect a number of records which may not be material individually but could have a combined material impact;
- static and reference data items which are spread throughout the internal model (eg exchange rates, correlation matrices and organisational hierarchy);
- data items that are key to stress scenarios (eg any error in these may be compounded under the stress scenarios and therefore may have a higher impact); and
- instances where data error made an item appear immaterial when it was material.

4.20 The Act III process demonstrated that many firms had not reviewed data flows to identify areas where there was a greater probability of a material data error and to verify if their existing control framework was sufficient to mitigate the risk. It is useful to start by itemising the key data operations, identifying the key data fields required for these to work, then finally (and crucially) assessing the risk and impact of error and appropriate control requirements. The initial risk and impact assessment will be based on inherent risk. Controls are then mapped in and the residual risk re-assessed to see whether it falls within the firm's risk appetite.

4.21 For example, one firm proxied residential property values by taking the most recent known independent valuation, then projecting this value from a quarterly property index for the appropriate geographic region to the present date. This implied the following data operations:

- select property detail, including location code, last property value and valuation date;
- map property location code to index region;
- map valuation date and present date to quarterly date series used by index, in order to impute index value at those dates, perhaps using adjustment or interpolation, then impute index return between those dates; and
- impute current 'proxy' value to property at present date using imputed index return.

4.22 This raised the following questions:

- Has every property detail been selected from the database, ie is the property data complete in the sense that every record which should have been selected, has been selected?
- How accurate is the mapping from the location code used in the database, to the code used by the index?

- How frequently is the mapping checked to account for changes in code, or errors in the code originally entered?
- Is the bucketing method used to map the imputed index value to the valuation date fit for purpose? For example, if the valuation date and current date are rounded down to the latest quarterly dates, is this likely to introduce a systematic bias?
- Finally, is there some strong appropriateness monitoring control which reconciles realised transaction values to the proxy value at transaction date, perhaps encoding this by a regression coefficient? Such forms of ‘back-testing’ or other forms of end-to-end reconciliation are a powerful test. Whenever two data sets (here, realised values and proxied values) are reconciled in this way, it is helpful to represent them both on the diagram, preferably linking them by a control, with a cross-reference to a control description.

4.23 The 2012 report noted ‘firms who successfully conducted such a review’¹ generally relied on indicators such as complexity of the data transformation process, complexity of the controls, level of manual intervention, historic evidence of error, reasonableness assumptions about error, etc’ (2012.4.21).

Finding 5: data directory

4.24 The 2012 report noted that the underlying purpose of a ‘data directory’ is to ensure that there is good governance over data quality, and to ensure firms think carefully about their approach to governance². It also noted that the data directory is meant to be a documented repository where different users can go to understand which data is being used in the model, where it comes from (‘source’), how it is used (‘use’) and what its specific ‘characteristics’ are.

4.25 The data directory continued to be a challenge for many firms. If the data directory is too simplistic for future use or too complex to maintain, this will defeat its intended purpose. A common problem was that firms documented data items at an unnecessarily detailed level of granularity, so that the directory was almost as complex as the system itself. Many directories document data at the level of database fields or record types (eg INV_SWP_CUR_PAY), confusing a ‘data dictionary’ for use by technology, with a ‘data directory’ in the Solvency II sense, ie a repository where different end-users (not technologists) can go to understand how the data is being used in the model.

4.26 The Act III process showed clearly that firms were not able to represent data flow in an appropriate way, and in particular were not able to identify the location of material control points to enable the successful scope, design and execution of data audits. Many firms submitted data flow diagrams which had data from one risk factor joining or mixing with data from a different risk factor in a way that made no sense. Other problems identified were:

- Lack of narrative. A number of data flow diagrams were sent without any kind of accompanying description.
- No categorisation. It is helpful to categorise broadly the type of operation performed on the data, eg aggregation, transformation, cleaning, joining or enrichment. If there is a join between data sets, indicate the appropriateness checks that have been performed. Note that cleaning of data is a kind of transformation, which ideally will leave a trail that

¹ That is, the impact and risk assessment of the data.
² See 2012.4.5 and Delegated Regulation 2015/35 Article 244(g).

includes the original raw data, the data used to inform the cleaning process and the cleaned output. It is also helpful to indicate where the check is 'correspondence' or 'consistency'. A correspondence check compares a sample of data received downstream, to data sent from a reliable source upstream. A consistency check uses known properties of the data (maxima and minima, averages, properties of the distribution) to detect possible errors. Most so-called data profiling is a form of consistency check.

- Circular or empty descriptions. This was a common weakness, applying to all the submissions to some extent. A circular description is when an operation is described in a way that adds little or no information about what the operation is. For example 'data is validated to ensure it complies with standards' is no more than a statement of what 'validation' means. It does not summarise what the validation process is, its purpose, and what it validates. An empty description is one that states what would be expected to be known already. For example 'each table is designed to store particular data.'. What kind of data? Where does the data come from? Where does it go? 'The manual download process is subject to manual controls to ensure the data is complete and accurate'. What counts as complete and accurate?
- Failure to identify key assumptions. There may be hundreds of assumptions underlying a model run or transformation process. However, only a handful of fields are likely to be economically relevant. It is generally helpful to identify these and place them in the narrative. For stresses, define the key risk factors used. For reconciliations, describe what is being reconciled (record count, balances, tracer IDs etc), and the same applies to calibration. What is calibrated to what? For example, is it a mathematical distribution to an empirical one? Is it the output of a proxy model to the corresponding 'full' revaluation?
- Not 'joined up'. Joins between data sets are always important. The risk of a failed join is usually high. The join key¹ nearly always requires mapping (system X represents geocodes² in one way, system Y in another). Poorly designed joins can lead to dropped records particularly when implemented in spreadsheets using the vlookup function. Nearly all firms have a problem representing or recognising the existence of joins. Every balance sheet valuation requires joining exposure data (terms and conditions, model points) to valuation data (prices, rates, qx³ tables). Every capital system requires joining the same exposure data to distributions or historical series of valuation records. A join nearly always requires some test for appropriateness of data, for example when a portfolio of single stocks is joined to historical index data. Every time the PRA reviewed a flow diagram, it looked for such joins, but they were rarely included. Every form of enrichment (for example, when bond ISINs are enriched with static bond reference data such as coupon details) also requires a join between the enriching and the enriched data, yet many narratives refer to the enrichment of data, without referring also to the two data sources which are necessary to do this. Identify the attributes which are used to join data (eg age and gender used to join policy to longevity data; and postcodes used to join property exposure data to catastrophe model event footprint).
- Other missing operations. Often the PRA received data flow representations where it was clear that a certain operation must exist, yet none was described. For example, if the model input referred to sensitivities or to spreads, this implied that bond sensitivities must have been calculated at some point. One firm referred to the 'de-notched' ratings used

¹ That is, a common cross-reference which ties one set of data to another.

² A geocode is any kind of code that represents a place. It could be a postcode, a longitude/latitude reference, it could be a proprietary method used by a system vendor.

³ qx is mortality death rate.

within the model. This implied: i) a join between the de-notched exposures and the spreads; and ii) a de-notching process requiring some form of control, but where neither the process nor the control was explicitly referred to. Gridding or bucketing was another process which must exist, but was rarely mentioned. For example, in order to join longevity data to annuity cash flows, the data must be structured and joined at cash-flow level, implying a process to generate the cash flows to predetermined periods. The example data flow diagrams in the appendices use subscripts such as 'm' or 'p' to indicate where such gridding might exist. This can also be used whenever there is an aggregation of data, such as when monthly cash-flows are aggregated to annual, or when policyholder data is aggregated by age.

- No control description. If there is no control description attached to the description of a data operation ask why no control is needed, or whether there are actions to address. If an operation really has no control, ask if it is important enough to describe in a high level data flow.

4.27 The following table lists some real examples from narratives provided by firms, with suggested more helpful alternatives.

Real examples	More helpful
'Files are extracted from market data systems and uploaded into the asset data warehouse via a data-stage routine.'	'Current interest rate curves, foreign exchange and credit spread data curves are extracted from the market data system and uploaded to the data warehouse.'
'Data is validated to ensure it complies with standards.'	'Exposure data is checked for completeness of records, focusing on postcode and construction details.'
'The XYZ system is used as a repository and processing function to consolidate the source data received from the Asset Data Warehouse.'	'The XYZ system consolidates policy data from the warehouse by aggregating into age groups.'
'A variety of tables store fund, asset and valuation information which are then reconstructed into a series of files in a specific format that can be input into the model.'	'Fund and asset valuation details are extracted from the warehouse and transformed into the data model expected by the model system.' (If there is enrichment or change to the data, please specify the enriching data source.)
'The requirements for creating Model Point Files have been specified in a series of Requirements documents which have been tested as part of the development of the process and reviewed through the Change Control Process.'	Not necessary. Where processes are entirely automated, it can be assumed that errors will be picked up by the general computer controls part of the audit. Any manual parts of model point creation should be identified and assessed for impact and risk (complexity). Controls for risk should be clearly explained.
'Reports are produced to assess whether the controls are operated, sufficient and complying with regulation.'	'Completeness and accuracy controls (checksum and average spread) are operated for each model run, with an appropriate exception and escalation process.'
'Controls validate that the assumption data required for the Model has been obtained.'	'On each run, assumption data held in user defined tables is checked against a template to ensure consistency. Changes to the template are journalised, and differences signed off by head of capital modelling.'
'The initial input of data is subject to inbuilt systems controls.'	'Exposure data is input manually, subject to a number of preventive checks including postcode and construction details. Omissions must be authorised and recorded by exposure manager. Detective checks, including average room height and median sum assured, are provided to the capital model.'

Finding 6: data classification

4.28 The 2012 report noted: '[C]lassification is generally only useful if it corresponds to a common risk, or impact, or control method or other characteristic relevant to data governance. For example, asset data and policy data (as well as reinsurance data) represent underlying contracts with third parties. Accuracy and completeness of the contracts may

require matching paper records of terms and conditions with electronic ones, and a method of records retention that reflect the sometimes long-lived and stable nature of a contract. By contrast, observational data such as credit ratings data, or lapse rates, are more volatile, and may require frequent monitoring and control' (2012.4.29).

4.29 The point and purpose of data classification were still not well understood by many firms. Reasons for doing this include the:

- need to understand how a capital model works without understanding the specifics of the data it uses, which at its simplest involves joining of exposure data to risk data to produce a financial distribution from which a 'Value at Risk' can be estimated;
- ease of identifying where data is missing on the representation. For example, if input data consists of bond terms and conditions, but output data consists of asset valuations, this suggests that pricing or rates data have been joined, and the representation is missing something. This will also indicate where a control for appropriateness may be required. Similarly, if input data is base improvement, but output data is a distribution of possible improvements, this suggests a distribution or stress of improvement factors has been used; and
- overall understanding of controls required. At some point in the flow, exposure data will have to be joined to valuation data (for base balance sheet) and to risk data (for capital requirement). One firm showed a narrative where policy data was input, and interest rate valuation was the output. Appropriate classification of data would have enabled the firm to question the logic of this narrative. In another case, exposure data for multiple locations was joined to valuation data for a single location, resulting in a substantial capital error. Classifying exposure data correctly would have avoided this.

4.30 2012.4.30 stated 'referential data (such as bond coupon definitions, global policy definitions) or configuration data (such as model parameters, and assumptions) determines the meaning or value of potentially large numbers of records, and therefore requires careful control against error, such as access and version control'. There was still not enough focus on the effect of assumptions, such as static and reference data items which have a wide impact across the internal model (eg exchange rates, correlation matrices and organisational hierarchy).

Sub-risk 4: Controls over data quality

Demonstrating the effective operation of data quality controls by evidencing and articulating the controls, and by showing how the control process is operating consistently in business as usual.

Finding 7: data quality controls

4.31 The 2012 report noted that nearly all firms had difficulty in demonstrating the effective operation of data quality controls, and that this was primarily due to a lack of evidence of controls (2012.4.31). It commented (2012.4.34) that firms must be able to articulate the nature of data quality controls, and to demonstrate how the process is operating consistently with appropriate controls in business as usual.

4.32 This was still problematic. For the Act III process firms were asked to provide (for sample risk factors) an end-to-end flow from source systems to internal model. Many struggled to provide an appropriate narrative covering the types of operation on data, and found it difficult

to identify the risk and impact, and the required control characteristics of the operations. Some firms had an informal control process, but there was a difference between evidencing a control (such as by sign-off, or a log) and documenting it (a 'how to' document explaining how to operate the control). The process of documenting a control nearly always suggests ways of improving control design, and without such documentation it is hard to identify situations where escalation is appropriate.

4.33 It is helpful, as noted in 2012.4.32, to use the data flow diagram to determine the key data operation points, in order to match controls to those points to ensure that appropriate checks have been completed.

4.34 Although many firms used data profiling (2012.4.33), subsequent reviews suggested some firms were not. Downstream functions therefore, run the risk of relying on data which is incomplete and inaccurate, possibly inappropriate.

Finding 8: third-party data

4.35 Many firms were still relying on controls operated by a third party despite having no mechanism to obtain assurance over the control environment and without independently validating the external data received (2012.4.36). For asset management, suitable controls would be straightforward to implement by, for example, sampling alternative third-party data, but be cautious of 'false reconciliation' where both data sources are reliant on an erroneous source further upstream. One firm relied on sourcing bond reference data to check the terms and conditions of its asset portfolio. However, the upstream asset portfolio relied on the same external source which sometimes contained errors. The correspondence control did not fail, even though there were errors in the data. Large or concentrated positions present a significant risk, but terms and conditions can always be checked against the original legal documentation on a sample basis. The same consideration applies to policies or other contracts, both retail and wholesale (eg reinsurance).

Sub-risk 5: IT environment, technology and tools

Spreadsheets and other user-developed applications are a form of information technology, and all information technology needs to be appropriately controlled.

Finding 9: end-user computing (EUC)

4.36 Spreadsheets and other end-user applications (2012.4.39) remained common in capital and balance sheet modelling. The PRA does not have a view on whether end-user computing (EUC) is appropriate, as it is a form of IT, and all IT needs to be appropriately controlled. Where EUC is material to the internal model data flow, the PRA will be looking for appropriate controls for data quality such as reasonableness checks, input validations, peer reviews, systems environment configuration, logical access management, ongoing change controls (development, build, systems and user acceptance testing) and release management (including implementation and operational testing), disaster recovery, and documentation.

4.37 Automation of spreadsheets reduces the risk of manual error (2012.4.42), but can introduce different problems such as reduced oversight, inadequate transparency about the extent of linking and the proliferation of nested spreadsheets and the attendant issue of 'broken links'.

4.38 The 2012 report did not engage comprehensively with cyber risk. This is likely to be an area of increasing focus, following alerts and increasing concerns about security as firms move away from localised application and onto networked platforms. As noted in the Bank of

England Financial Stability Report,¹ cyber attacks can threaten financial stability by disrupting the provision of critical functions from the financial system to the real economy. The Financial Policy Committee has recommended that resilience testing be a regular part of core firms' cyber resilience assessment. Insurers providing cover for cyber or business interruption are also indirectly exposed to cyber risk.

Finding 10: IT infrastructure

4.39 Complex IT implementations (2012.4.44) can be challenging to manage without a clear definition of user requirements, design, testing and appropriate controls for effective operation in business as usual. This continued to be an area of risk. One firm took seven years to implement a tactical system, and still has no strategic system for its upstream administration processes.

5 Maintenance regime

5.1 By 2016 some firms will have gone through IMAP and received model approval, some may be in the process of applying for model approval. Approval is not the end of the process, however. Solvency II Directive Article 36 specifies that supervisory authorities shall review and evaluate the strategies, processes and reporting procedures used by firms to comply with the Directive, referencing the requirements for technical provisions (Solvency II Directive chapter VI sec. 2) and capital calculation (chapter VI secs. 4 and 5). These include the requirements on data quality; firms may wish to consider how they will evidence these as part of the reviews conducted by authorities.

5.2 The PRA data audit is an established means of providing assurance to the PRA. Firms may wish to consider the frequency at which they conduct such reviews. They may also wish to consider that the Act III review process (the purpose of which is to establish whether the PRA can place reliance on the data audit) may form part of the PRA's data quality assurance programme. This document has explained this process, and has outlined the challenges some firms have found in satisfying it. The appendices outline possible ways of evidencing data flow to the PRA, by means of data flow diagrams, annotated control points and narrative description.

5.3 Auditors (internal or external) may wish to consider whether the information provided by line 1 or line 2 is sufficient to perform an audit which is adequate in terms of scope, design and execution.

1 Part A 'Cyber risk', July 2015; www.bankofengland.co.uk/publications/Pages/fsr/2015/jul.aspx.

Appendices

-
- | | |
|----------|-------------------------|
| 1 | Credit data flow |
|----------|-------------------------|
-
- | | |
|----------|-------------------------------------|
| 2 | Catastrophe (CAT) model flow |
|----------|-------------------------------------|
-

Appendix 1: Credit data flow

1.1 The diagrams and descriptions below represent one of many possible ways of representing the data flow from exposure, valuation and risk data systems to balance sheet and internal model. The method of estimating capital requirement described is also one possible estimation methodology, and is not general guidance under the Financial Services and Markets Act 2000. The flow only looks at a few ‘modelling assumptions’, ie interest rates and credit spreads. There are many more, eg exchange rates (payments in foreign currencies), retail prices index assumptions (index linked bonds), modified durations (most commonly used in the modelling of futures), Libor (floating rate notes) and volatilities where embedded options exist (callable and puttable bonds). Each of these needs to be ‘joined’ to the modelled data.

1.2 Operation [1].¹ Receive m^2 bond prices to value the bonds, identified by using a unique reference such as ISIN. Depending on data source (fund managers may provide in one file, data vendors most likely in separate files or downloads) price data may arrive in the same file or table as underlying bond exposure data and reference data, but it is still helpful to keep this separate in the diagram in order to distinguish data used for valuation from data used to define exposure.

1.3 If the data already includes the actual valuations this may be included in the narrative.

1.4 Possible controls [1]. If prices are used, controls can include consistency checks (looking for known properties of the data, including expected daily changes, consistency with bond spreads, yields etc) and correspondence checks (using other data suppliers). For correspondence checks be careful to avoid false reconciliation, ie using an apparently different source which really derives from the same ultimate source upstream.

1.5 If unlisted assets are included in the flow, appropriate controls may be recorded.

1.6 Operation [2]. Prices are joined to bond exposures, either by bond identifier or by record, to value the bonds. Bonds without prices can be branched to a separate process for alternative valuation. Ensuring the join links all input data and reference data together presents challenges, eg multiple identifiers for the same bond from different sources: identifier licensing issues present challenges to the client requesting the data.

1.7 Possible controls [2]. If bonds do not have a price, consider an exception report identifying the number of bonds affected, plus the value assigned by the default valuation. Operation [3]. Join bond unique references (eg ISIN) to issuance reference data to create the full terms and conditions (T&Cs) of the bonds.

1.8 A row from the answer table output by [3] may look like the following:

¹ References in square brackets, [1] etc relate to those in Figure 1 of this appendix.

² The subscripts (m, n, i, j etc) represent in all cases the dimensions of the data, ie the number of records for each dimension required.

Name	ABC BANK 06/46 MTN
Issuer	ABC Bank PLC
ISIN	XS0247840969
Coupon	4.75
Payment count	Annual
Coupon date	24 Mar
Issue Date	24-Mar-06
Issue Volume	600.00 m.
Issue Currency	GBP
Maturity Date	24-Mar-46
Bond Denomination Currency	GBP
Settlement Currency	GBP
Subordinated	Yes
Last price	90.74
Last price timestamp	07/06/2013 02:13

1.9 These legal entity identifier (LEI)¹ details are received, either from a separate source, or from the original data source (ie [1] above), and joined to the m exposure records via an appropriate join key, eg bond identifier code (ISIN), or record row.

1.10 Possible controls [3]. The join is typically a mechanical one, which according to the PRA's reviews is rarely checked. If sourcing the referential data externally in order to check the T&Cs received from the administrator, be careful of false reconciliation, where the administrator has used the same external source. The primary data for bond T&Cs is the legal term sheet or 'tombstone', which is typically a free-form document that must be manually checked against the electronic representation. T&Cs are frequently 'shoe-horned' into electronic records, which can give rise to material valuation errors. To mitigate this risk, consider periodic risk-based sampling to compare T&Cs in upstream systems of record, with those in downstream systems. If there are unlisted assets with no T&Cs, consider using exception reports to identify the percentage of such assets, both by number and by value.

1.11 The join is typically parent (Primary Key of Issuer name ie 'ABC Bank PLC 'T&C for bonds), to child (foreign key of 'Bond issued by 'ABC Bank PLC') which the PRA's reviews suggest is rarely checked. If sourcing the parent referential data externally in order to check the T&Cs received from the administrator, be careful again of false reconciliation traps.

1.12 Operation [3a]. Enrich m records of T&Cs with p interest rate cash-flows to create m*p data records.

1.13 If the source used to grid into p maturities is not the same as the one used firm-wide, it is helpful to indicate this. At some point the number of grid points (currently p) will need to agree with the number of cash-flow points used by the rest of the firm's systems. Each cash-flow will have an interest rate applied to it from a range (number of years to maturity) in the sequence given by the yield curve.

1.14 A row from the answer table output from [3a] could look like the following:

¹ See [www.gfma.org/initiatives/legal-entity-identifier-\(lei\)/legal-entity-identifier-\(lei\)/](http://www.gfma.org/initiatives/legal-entity-identifier-(lei)/legal-entity-identifier-(lei)/).

Name	ABC BANK 06/46 MTN
Issuer	ABC Bank PLC
ISIN	XS0247840969
Coupon	4.75
Payment count	Annual
Coupon date	24 Mar
Issue Date	24-Mar-06
Issue Volume	600.00 m.
Issue Currency	GBP
Maturity Date	24-Mar-46
Bond Denomination Currency	GBP
Settlement Currency	GBP
Subordinated	Yes
Last price	90.74
Last price timestamp	07/06/2013 02:13
Sensitivity 6m	£3m
Sensitivity 1y	£3m
Sensitivity 2y	£2.75m
Sensitivity 3y	£2.5m
Sensitivity 4y	£2.2m
Sensitivity 5y	£2.1m
Sensitivity 7y	£2m
Sensitivity 10y	£2m
Sensitivity 15y	£2m
Sensitivity 20y	£2m
Sensitivity 30y	£2m
Recovery	40%

1.15 Possible controls [3a]. If firm-wide data are used, it is sufficient to inherit the firm-wide system of controls over accuracy and completeness. Note any materially incomplete sets (eg significant missing maturity points in the range p expected). At the point where the interest data are joined to cash-flows there should be a test of appropriateness. For example, if interest rate data is available only up to 15 years, but asset cash-flows are available for 50 years, consider some appropriateness test that the interest rate data is good for purpose (for example see Commission Delegated Regulation Article 231 (1b) which requires that the data are recorded in a timely manner and consistently over time).

1.16 Operation [3b]. Bond T&Cs are used to project cash-flows for p maturity dates, joined to interest rate data to compute, for example, present values.

1.17 Possible controls [3b]. Note, as above, the risk that the maturity points assumed in the cash-flow projections may not match the maturity points of the interest rates (if, for example, an extra point has been added to the firm-wide data and the projection system cannot recognise this because the points are hard-wired). Good practice, as with any join between data where the same dimensions are expected, is to define cash-flow points in user tables, and to have an exception process identifying dimensional changes required to match incoming data.

1.18 As with any heavy to light process (creation of model points, replication formulae, sensitivity calculation) there may be a periodic process which perturbs both the heavy and light model by the same change in risk factor, comparing the results on an x-y chart, or performing regression analysis. Commission Delegated Regulation Article 240(2) requires that

the specification of profit and loss shall be consistent with the increase and decrease of the monetary amount underlying the probability distribution forecast, ie consider whether changes in valuation given by the heavy model are broadly consistent with changes implied by the light process used by the capital model.

1.19 Operation [4a]. Issuer ratings are grouped using interpolation; ratings data on issuers is received from a third-party source. Agencies have specific formats, eg rating for ABC Bank PLC is typically of the format AA minus or B minus. Ratings are 'de-notched' into a set of values mapping onto the same structure to produce a range of values consistent with those used by the historical time series (see paragraphs 1.25-1.26 below).

1.20 Possible controls [4a]. Consistency of ratings by issuer: the ratings data will have dimension (ie total number of issuers), which must match the corresponding dimensions of the historical data. If a new issuer is added to the historical data, this must be added to the issuer data, so consider a control to identify this by exception. Is there a similar process to identify any new rating type in the maintenance of historical data? Missing issuer ratings can be shown in exception reports.

1.21 Operation [4b]. Ratings data by issuer is joined to sensitivity data, joining on issuer, mapping where necessary where issuer descriptions are different. Maturity data is typically aggregated at this point, either by rating or by issuer depending on whether the risk data is rating or issuer specific to give sensitivity to loss of the issuer as sensitivity to changes in interest rates.

1.22 An answer table row from [4b] may look like this:

Issuer	ABC BANK PLC
S & P	AA-
Moody's	A2
Fitch	0
Denotched	AA
Sensitivity 6m	£31m
Sensitivity 1y	£31.5m
Sensitivity 2y	£32m
Sensitivity 3y	£33m
Sensitivity 4y	£34m
Sensitivity 5y	£35m
Sensitivity 7y	£36m
Sensitivity 10y	£36m
Sensitivity 20y	£35m
Sensitivity 30y	£31m
Recovery	40%

1.23 Possible controls [4b]: There is no conventional format for issuers, thus the same issuer can be represented in different ways. This problem is usually resolved by manually matching the fields in the first instance to create a mapping table, then running the mapping automatically for each model run. New issuers, or change to issuer name, will create exceptions which need to be handled appropriately, and consider controls to ensure that records are not dropped by a faulty join.

1.24 Operation [5]: Data cleansing and refreshing of the historical ratings data and the attributes of spread and default rate etc. This is updated or refreshed, and cleaned if necessary.

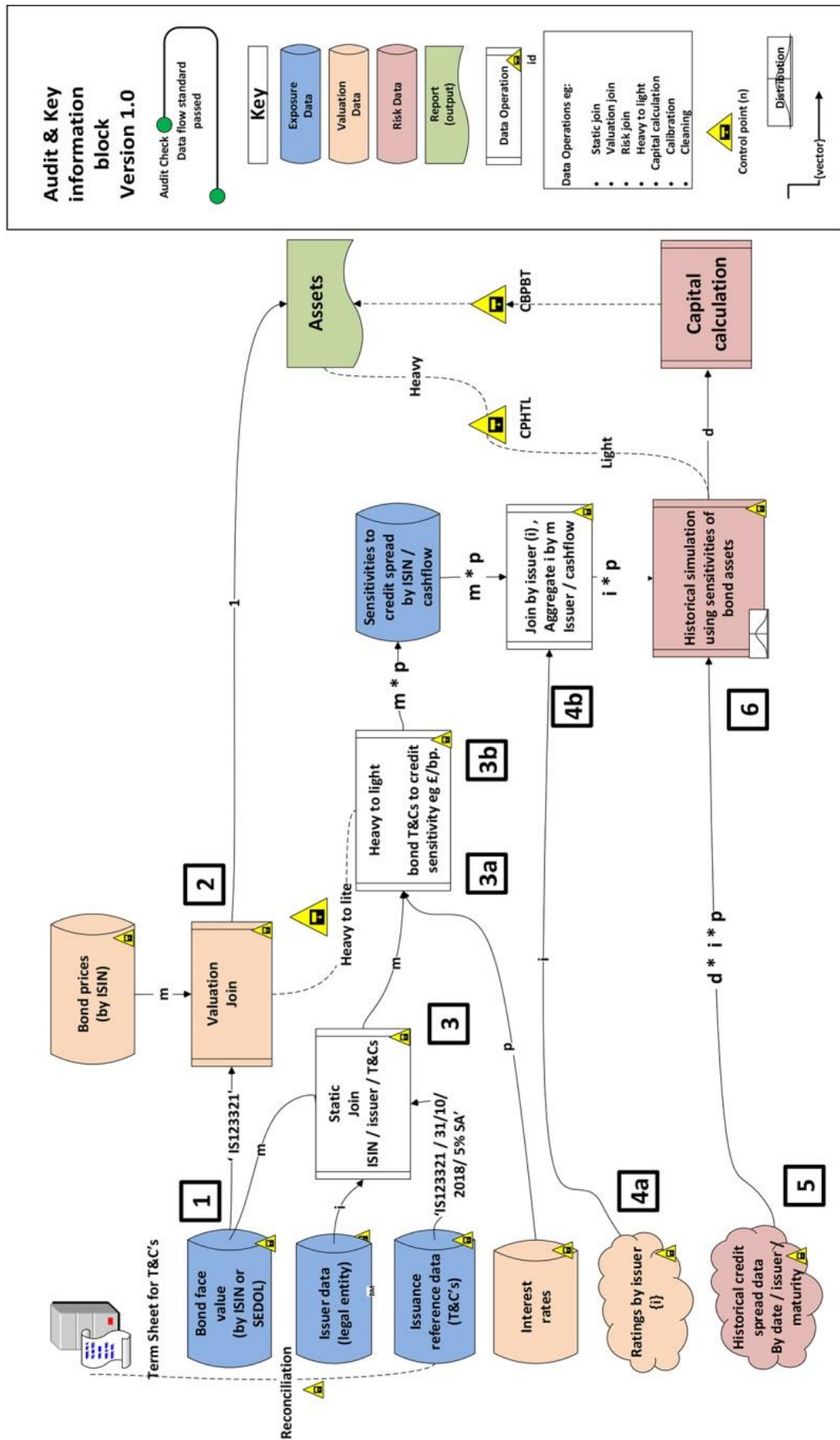
1.25 Possible controls [5]: Good practice is that any cleaning of data must be journalised, ie record changes by user and time stamped, in a way that is reversible. Reasons for the change may be identified. There should be access controls which reflect segregation of duties (doer, checker).

1.26 The historical file is typically large, being 'l * p * d', where l is the number of issuers, p the number of 'maturities' and d is the number of days history. Alternatively, the data may only reflect spread level by rating band, and so will have dimensions r, where r is the number of ratings. The use of historical time series has been assumed. A commonly used alternative is the economic scenario generator or ESG. Here, d will be the number of simulations, with p and r as before.

1.27 Operation [6]: Historic spread data is joined to sensitivity data produced in [4b] to produce a distribution of n profits and losses, hence a capital requirement number at a specified confidence interval and time horizon.

1.28 Possible controls [6]: The periodicity of the historical data should be consistent with the time horizon assumed by the capital model. Solvency II models assume a time horizon of one year, hence the return period of the historical or simulated data should be consistent with this. This should be checked as part of the process whenever the historical time series (or the ESG) is changed. If inconsistent return periods are used, the appropriateness of this should be justified.

Figure 1 Possible data flow for credit default risk



Appendix 2: Catastrophe (CAT) model flow

1.1 Operation [1]¹. Exposure data is extracted from the policy system for all live policies via a series of batch jobs, merged with bordereaux exposure ('delegated authority') data held on spreadsheets, and uploaded into the data warehouse where it is checked for completeness and accuracy.

1.2 Possible controls [1]. For batch processes, were the original scripts adequately tested to ensure that data is complete, ie is the data expected the same as data received? Once tested, is there a formal change control process, and are changes in line with the internal model change policy? Is the functionality of the script broadly intelligible to the ordinary user? Is there a suite of reports that allow end-users to test for completeness?

1.3 Spreadsheets are typically used for receiving information from delegated authorities. The PRA has no view on the use of spreadsheets, which are a form of IT system, and which (like any IT system) must be appropriately controlled. Spreadsheet controls might include adequate testing for the process extracting data from spreadsheets, and a formal change control process just as for corporate IT systems. If, as is good practice, a firm has an EUC policy, it is encouraged to apply it. As well as generic guidance (access control, good spreadsheet organisation and design principles) the policy should ensure a focus on what the spreadsheet does, and particularly on spreadsheets that are part of the operational fabric of balance sheet and capital modelling, not just those close to the kernel of the capital model.

1.4 Are there controls to ensure that the exposure data is materially accurate? Controls might include consistency checks such as data profiling of material fields. The principle of proportionality suggests that data should be tested in proportion to the impact of potential errors on the capital calculation, on key risk calculations, and on the base balance sheet.

1.5 Exposure data from delegated authorities is inherently at risk of being incomplete and inaccurate. Firms should have an appropriate risk appetite for poor data quality, and should be able to quantify their data quality risk according to that appetite.

1.6 Operation [2]. Where any primary characteristics (eg construction, occupancy, year built etc) are missing, policy data is enriched using static assumptions. For example, construction detail is enriched using 'industry data' that joins construction data to geocode. Sums insured and construction data, aggregated by peril and geocode, are uploaded into the hazard module of the CAT model. Policy conditions are uploaded into the financial data module.

1.7 Some primary characteristics are adjusted so that the simulated losses from the model are consistent with the firm's own claim experience.

1.8 Where a system of defaults is used (eg unknown property type defaulting to 'house', unknown year built to '1955', unknown number of bedrooms defaults to 3 for 'house') the default methodology should periodically be assessed and challenged for prudence.

1.9 Possible controls [2]. Enrichment is a form of join, and all joins are a complex and error prone process. There should be controls to ensure that the mapping from the geocodes used by the CAT model to those used by the exposure system are credible, and that the risk of data

¹ References in square brackets, [1] etc, relate to those in Figure 1 in this appendix.

(record or field) loss through failed joins has an appropriate exception and escalation process. One firm assumed that an exposure spread across multiple locations was in a single location, causing a breach of its capital requirement. Another firm mistakenly mapped a significant number of policies to a building that was never insured.

1.10 Operation [3]. The CAT model is run by a third party using the exposure and policy T&Cs to produce event loss data.

1.11 Possible controls [3]. Most CAT models are proprietary 'black boxes' leased by the vendor to end-users. The business model of the model vendors is based on the quality of the risk data used in the model. This data includes geo hazard information based on the likelihood of a peril occurring in a particular place, which is simulated (dimension 'n' in Figure 1) to produce an event footprint, then a local intensity footprint of locations near to, and hence potentially affected by, the event. Data also includes 'vulnerability curves' which are a form of light model mapping intensity of risk factor (wind speed, ground movement etc) to percent of sum insured ('damage ratio').

1.12 As the model and data are proprietary, there are no controls that the firm can directly impose. However, the data used in external models is effectively a form of external data, and firms should note the provisions of Delegated Regulation Article 19(4), which effectively apply to the data output from the model. Include reconciliation of the sum assured, also reconciliation of claims implied by the model to the firm's own claims experience. In cases where the firm adjusts the model on the assumption that its own claims experience is 'superior', this is a form of expert judgement that should be clearly evidenced, routinely challenged and periodically tested. For example, claims experience typically applies only to frequent and low severity losses, whereas catastrophe events are essentially high severity low frequency events, meaning that large losses are unlikely to appear in claims data. Thus any such adjustment impacting a modelled distribution, and particularly the tail of the distribution, should be carefully justified.

1.13 Operation [4]. Re-simulation of model output. The firm's own model uses the 'n' simulations from the vendor model, and resamples and rescales to produce a capital number (eg for Solvency II at the 99.5th percentile of the loss distribution).

1.14 Possible controls [4]. Check the number of simulations from the vendor model is consistent with the number expected. Check that the distribution is appropriate by comparing the estimated technical provisions using claims or other data is consistent with the median of the distribution. Another appropriateness test, albeit crude, is to compare profits and losses from previous periods with the modelled distribution.

1.15 Assumptions data used to configure the model run are reported (eg as a control file) detailing the assumptions used, including checks performed and sign-off.

Figure 1 Possible data flow for CAT loss model

