

**UNITED NATIONS STATISTICAL COMMISSION and
ECONOMIC COMMISSION FOR EUROPE
CONFERENCE OF EUROPEAN STATISTICIANS**

**EUROPEAN COMMISSION
STATISTICAL OFFICE OF THE
EUROPEAN COMMUNITIES
(EUROSTAT)**

**ORGANISATION FOR ECONOMIC
COOPERATION AND DEVELOPMENT
(OECD)
STATISTICS DIRECTORATE**

Joint UNECE/Eurostat/OECD work session on statistical metadata (METIS)
(Geneva, 9-11 February 2004)

Topic (i): Functions of metadata in statistical production

**METADATA SYSTEMS IN STATISTICAL PRODUCTION PROCESSES -
FOR WHICH PURPOSES ARE THEY NEEDED, AND HOW CAN THEY BEST BE ORGANISED?**

Invited Paper

Submitted by Statistics Sweden¹

I. INTRODUCTION

After an analysis of the purposes and contents of statistical metadata systems, this paper discusses how these systems can be implemented more efficiently and maintained in a sustainable way in statistical organisations.

II. The purposes and roles of metadata in statistical systems

Figure 1 gives a schematic picture of a statistical (production) system. There are two main loops in a statistical system, the design loop (reflected in the upper part of the picture) and the operation loop (reflected in the lower part of the picture).

The design loop includes

- planning, construction, and implementation processes
- evaluation processes

The operation loop includes

- execution and monitoring processes
- use processes

¹ Prepared by Bo Sundgren, bo.sundgren@scb.se.

STATISTICAL (PRODUCTION) SYSTEM
S

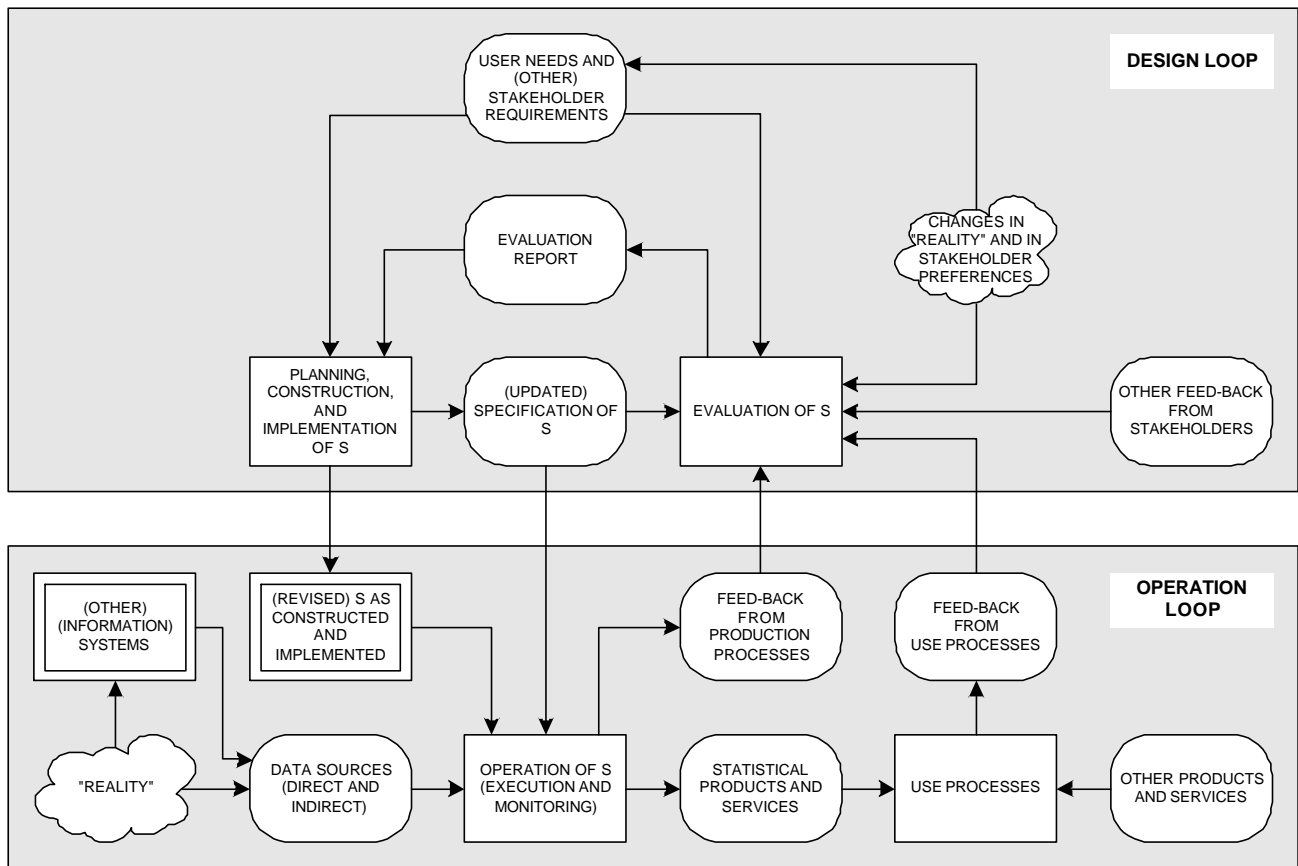


Figure 1. A statistical production system.

Most processes in a statistical system use and produce data, and they also use and produce metadata describing the data. There are numerous data/metadata sets and data/metadata flows in a statistical system. The data/metadata sets and flows are highly interrelated and they have to be so, both because they are logically dependent upon each other, and because – as we shall discuss later in this paper – statistical systems may become much more efficient and cost effective, if one systematically exploits the logical dependencies, e.g. by reusing and automatically transforming data and metadata as often as possible, thus avoiding costly manual operations like data capturing.

Main users of statistical metadata are users and producers of statistics. Other categories of stakeholders in connection with statistical metainformation systems are respondents, managers on different levels, and funders of statistical organisations. There are many categories of users of statistics with quite different needs and competence profiles: researchers, analysts, journalists, politicians, students, and the public at large, for example. Among the producers we may also distinguish between different roles, e.g. producer/planner and producer/operator. Very often the same person may perform a combination of these roles, so it may be better to talk about “usages” rather than “users” of metadata.

In short summary we may identify the following main kinds of metadata needs for the different stakeholder categories:

- Users of statistics need good metadata in order to identify, locate, retrieve, interpret, and analyse statistical data of relevance for their primary tasks. *Cf figure 2.*

- Producers/operators of statistical systems need metadata for the same purposes, but also for executing and monitoring the production processes properly, and for training new staff members. Cf *Figure 3* and (with more details for different kinds of operation processes) *Figure 4*, *Figure 5*, and *Figure 6*.
- Producers/planners of statistical systems need metadata for designing, constructing, and implementing statistical systems. Cf *Figure 7*.
- Respondents need metadata in order to understand why their participation in a survey is needed and for interpreting the meaning of the questions to be answered.
- Managers need metadata in order to evaluate different aspects of statistics production, including aspects of production efficiency, user satisfaction, and acceptance by respondents.
- Funders have similar needs as managers but on a more global level. They also need metadata that help them to balance the needs for statistical information against other needs that they have.

In all these cases of usages of metadata it is human beings who are the users of the metadata. In addition

- Software products and computerised data processing systems used in production and usage of statistics need metadata in order to function properly.

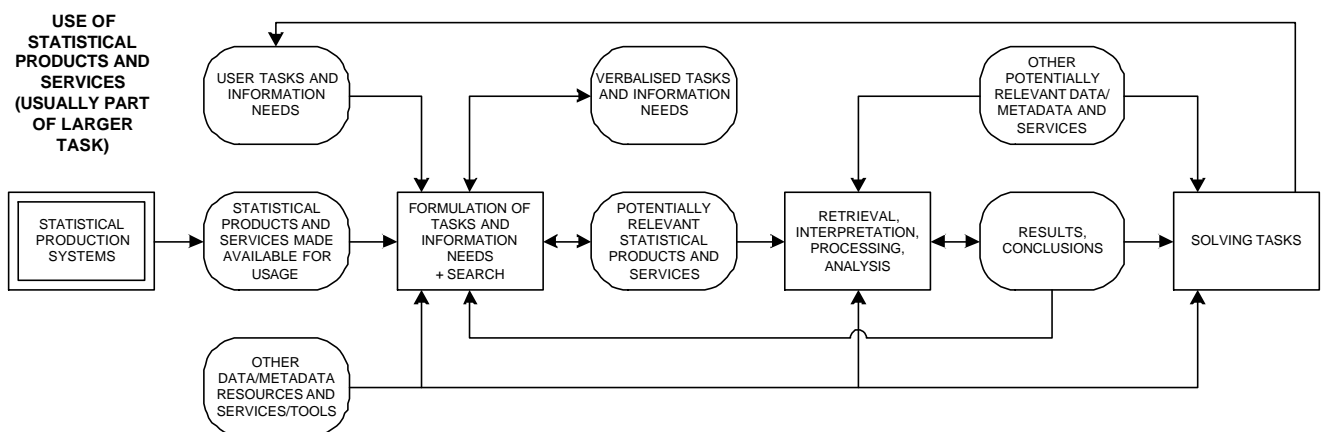


Figure 2. Usage of statistical products and services.

We may distinguish between metadata used by human beings and metadata used by software. Sometimes the same metadata, e.g. text labels, may be used both by human beings and by software, but usually there are quite different kinds of requirements on these two categories of metadata. Metadata used by software products usually need to be highly structured and formalised. Among human beings the requirements on metadata vary from user to user to some extent. Some users prefer more structured metadata than others. Professional users of statistical data may need more detailed and precise metadata, but on the other hand they usually have good background knowledge, which helps them to interpret data and metadata in a correct way. Casual users may not need metadata with the same precision, but on the other hand it is important for them that metadata are designed and presented in such a way that they interpret the data correctly, at least by and large, although their background knowledge may be rather superficial. Some users have plenty of time, whereas others have to digest data very quickly. Thus user requirements on metadata are different and sometimes even contradictory. This has to be taken into account when designing metadata and metadata systems.

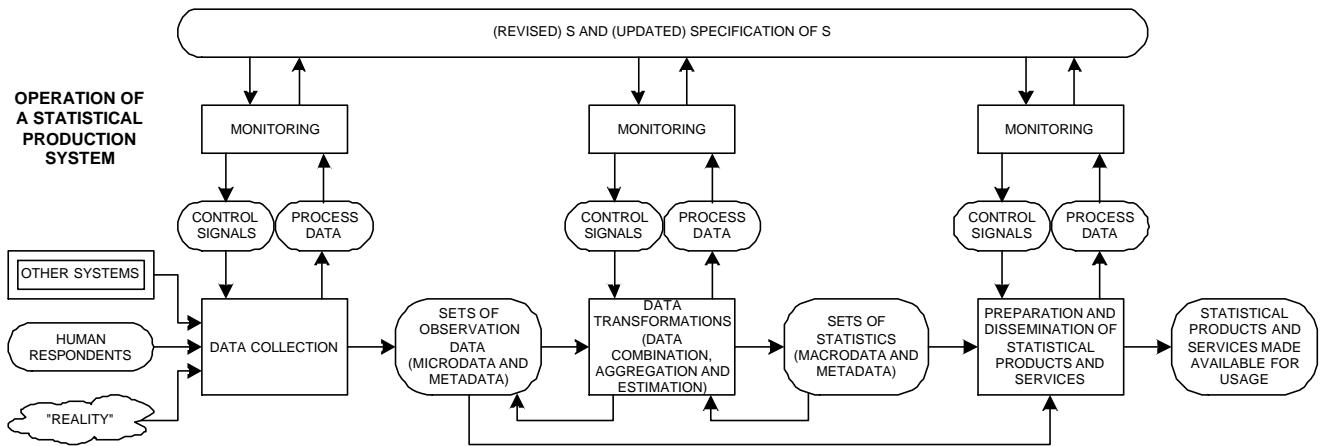


Figure 3. Execution and monitoring of the tasks of a statistical production system.

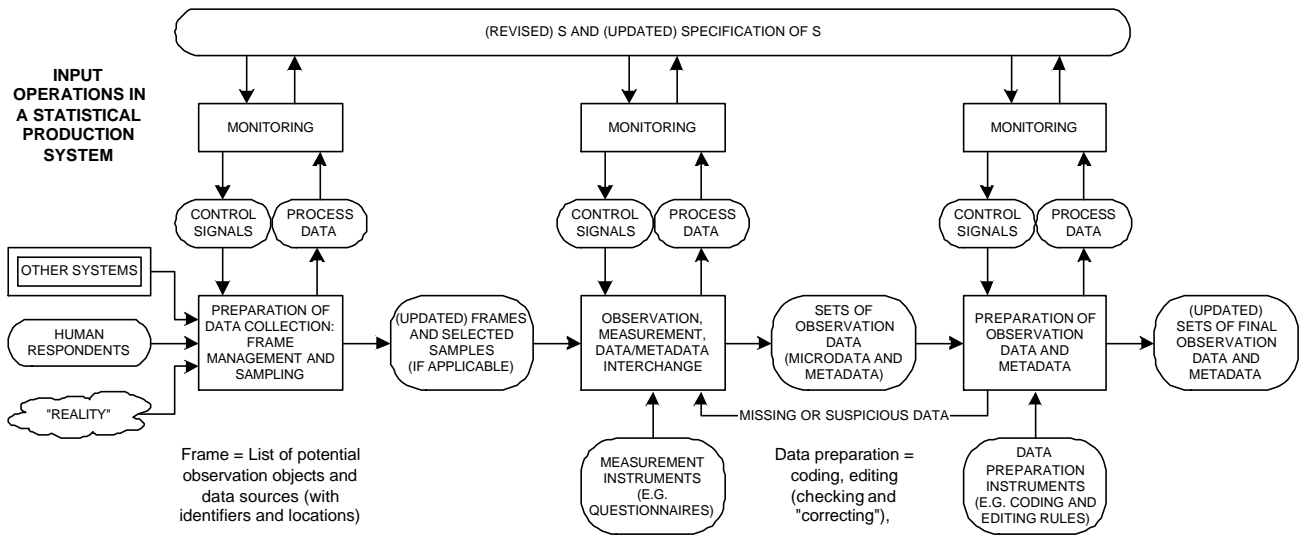


Figure 4. Input operations in a statistical production system.

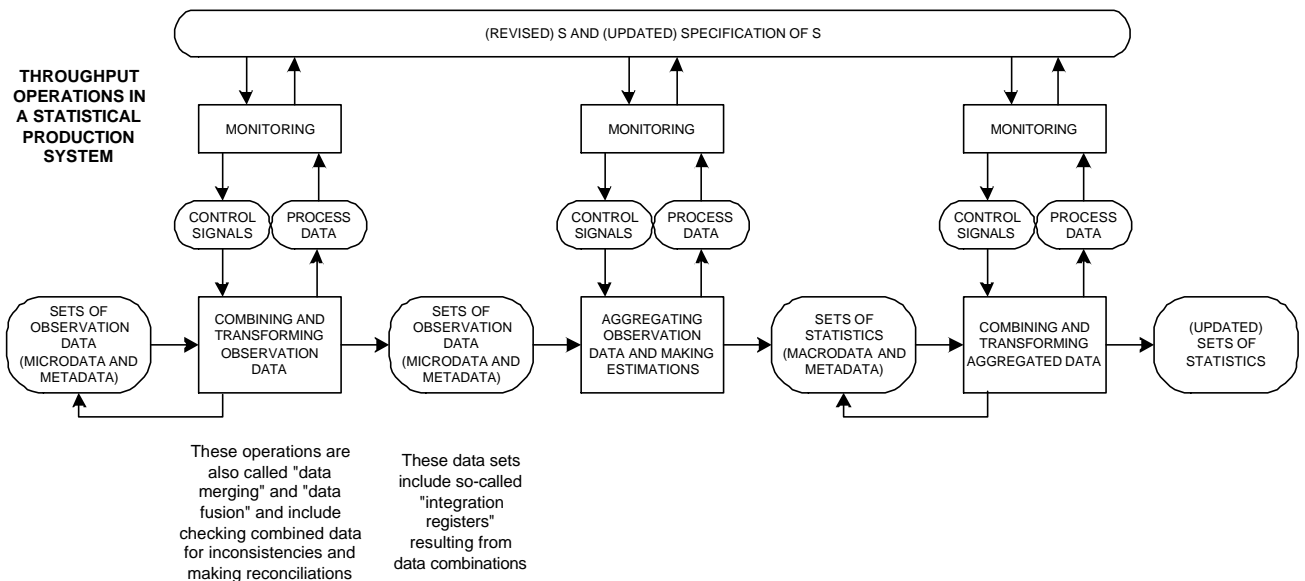


Figure 5. Throughput operations in a statistical production system.

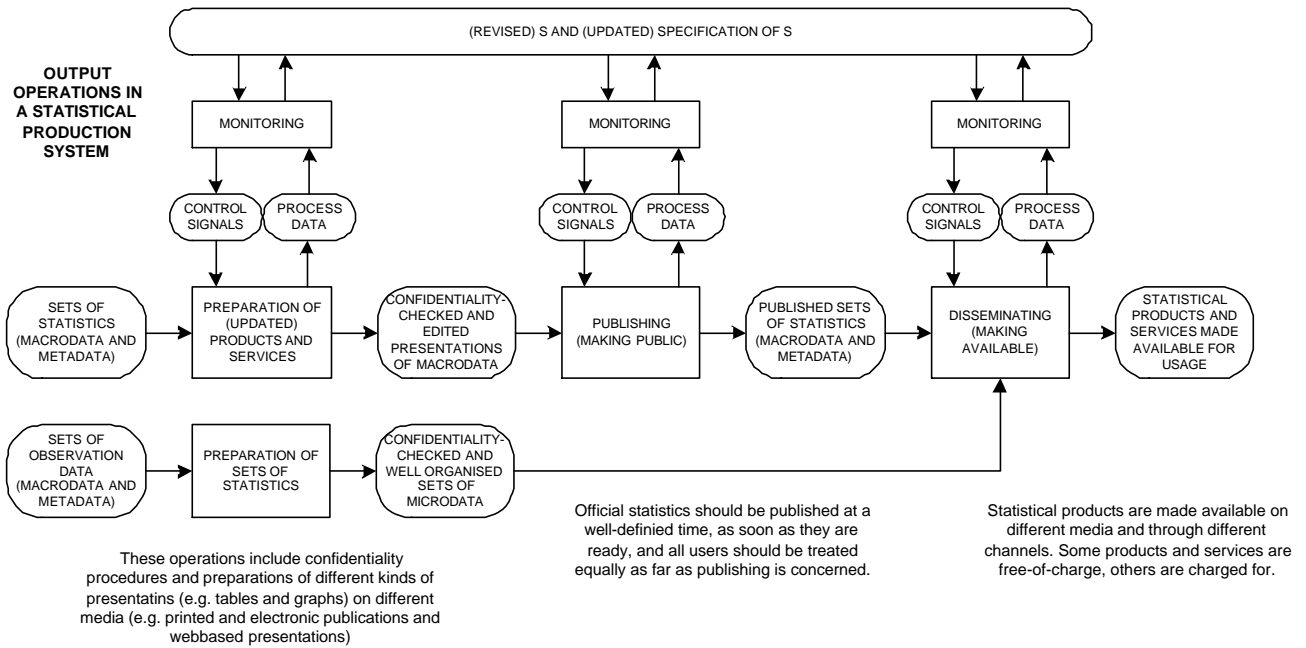


Figure 6. Output operations in a statistical production system.

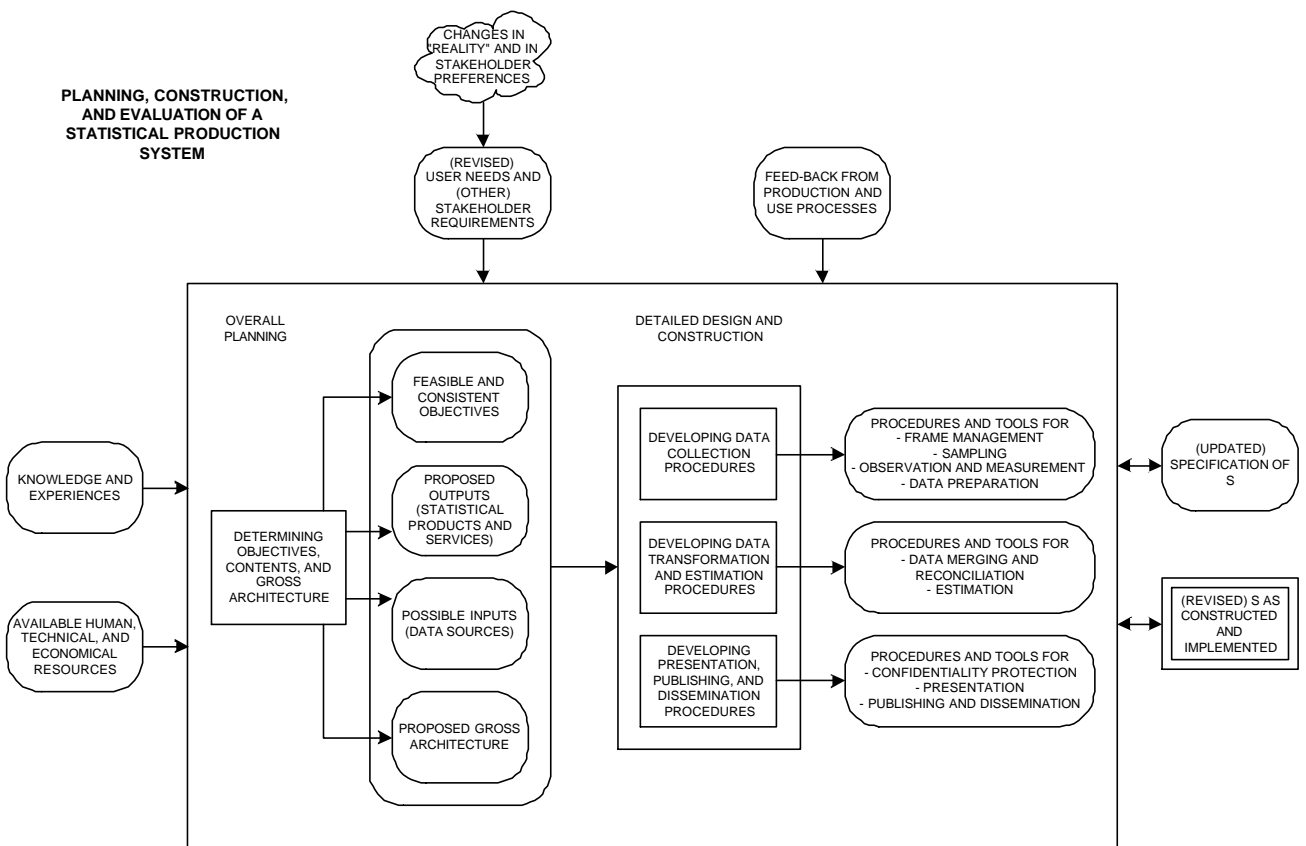


Figure 7. Planning, design, construction, and implementation of a statistical production system.

III. What should statistical metadata inform about?

All metadata in a statistical organisation can be organised around three main object types:

- data sets: observation data (microdata) and statistics (macrodata)
- processes: transforming data (and other inputs) into data (and other outputs)
- resources: instrumental in controlling and executing processes

A. Data sets: contents and quality

Since the ultimate purpose of statistical organisations and statistical systems is to provide users with statistical data (macrodata and microdata), data sets can be seen as the most fundamental object type; all statistical metadata can be seen as directly or indirectly linked to sets of statistical data.

Users are interested in the contents and quality of statistical data.

The contents of statistical data may be defined in terms of (for microdata) object characteristics, and (for macrodata) statistical characteristics.

An object characteristic is defined as

- the values of one or more variables (V)
- for the objects in a certain set of objects (O), also called population

A shorthand notation for a statistical characteristic is O.V. An object characteristic may be considered and described

- as originally conceptualised (the ideal characteristic or characteristic of interest)
- as operationalised after a design process (the target characteristic)
- as actually observed after the execution of a data collection process

A statistical characteristic is defined as

- a statistical measure (m) applied on an object characteristic (O.V)

A shorthand notation for a statistical characteristic is O.V.m. A statistical characteristic may be considered and described

- as originally conceptualised (the ideal characteristic or characteristic of interest)
- as operationalised after a design process (the target characteristic)
- as actually estimated

If we analyse the “deep structure” of statistical characteristics, we can identify at least the following (meta)object types:

- object type, population, subpopulation (stratum, domain of interest)
- variable, value set
- statistical measure

A population is defined by means of a population-defining property, e.g. “being a person living in Sweden”. Populations may be subdivided into subpopulations (domains of interest, strata); this is often done by cross-classifying the population by means of a number of variables, the classification variables.

Since both microdata and macrodata are associated with errors and uncertainties, statistical data are always estimated values, approximations of the true values. Statistics producers use quality declarations for describing

the discrepancies between estimated values and true values. Some types of discrepancies, e.g. sampling errors, are relatively easy to estimate quantitatively, whereas others can only be discussed verbally.

It should be noted that most quality aspects can only be given a partial description by the producer of statistical data. Whether certain statistical data are of sufficient quality for a certain intended purpose can only be finally judged by the user of the data. What is good quality for one purpose may be inadequate for another one. But the producer should aim at providing as good descriptions as possible of all aspects that may be important for a user trying to judge the usefulness of certain statistical data for a certain purpose.

Many statistical organisations, including Eurostat, have adopted standardised quality concepts, according to which the quality of statistical data is subdivided into a number of components like

- relevance: contents vs purpose
- accuracy: coverage, sampling, response, measurement, processing, model assumptions
- timeliness and punctuality: frequency, delays
- accessibility and clarity
- comparability: in time and space
- coherence: comparability and integratability of different statistical data
- completeness: the extent to which a certain domain is covered by statistics

B. Processes and resources

All kinds of processes in statistical systems are characterised by

- inputs to the process
- outputs from the process
- instrumental resources (for execution of the process and for control and evaluation of the process)

In an operation process the typical inputs (cf raw material) and outputs (cf intermediary and final products) are data and metadata.

Examples of instrumental resources for execution are:

- Data/metadata resources
 - forms, e.g.
 - data collection forms (questionnaires)
 - data storage forms (record layouts, database schemas)
 - presentation forms (reports, web layouts)
 - structured lists, e.g.
 - registers (authorised lists of objects of some kind)
 - classifications (authorised lists of values and codes, including relations between them)
 - thesauri (lists of terms, including relations between them)
 - rule sets: sets of editing rules (checks, imputations), sets of coding rules
- Intellectual resources
 - general knowledge: theories, methodologies
 - experiences: experience data from systems and processes
 - methods: e.g. estimation methods, algorithms, heuristics, procedures
 - models: e.g. data models, process models
- Human resources: organisation, staff
- Physical resources: buildings, technical resources, etc

Instrumental resources for control and evaluation are process specifications (typically in the form of some kind of documentation produced during the design phase) and process data generated by the process itself and fed back into the control mechanism. Process data inform about different aspects of the performance of the process, primarily such aspects that have a bearing upon the quality (e.g. accuracy, timeliness) and costs of the end-products. Process data for an operation can be obtained by monitoring and registering data about events that happen during the processing of individual transactions, e.g. during the capturing, coding, or editing of an observation.

IV. More efficient metadata management

The resources needed to develop, implement, operate, and maintain statistical metadata systems are considerable, both in terms of time, money, and efforts. Thus this work has to be carefully planned in such a way that it becomes efficient and accepted by all categories of stakeholders. The necessary planning is quite complex and requires active involvement and commitment of the top management of the statistical organisation.

One major challenge is to organise the basic collection and creation of metadata in an efficient way. Collection and creation of statistical metadata implies documenting different aspects of statistical processes and data. The people who are themselves actively involved in the production processes are obviously those who are best suited for this task, since they are in possession of the necessary knowledge. However, precisely because they already have this knowledge, they may not find it to be a task of high priority to document it.

Management could do several things here. They could show that they value the production of useful documentation and metadata by identifying this as a criterion for professional recognition and salary rises. Furthermore, they could discuss with the production staff, how metadata could possibly also be useful for the producers themselves, and give priority to metadata systems or system components that are able to serve both users and producers.

Capturing of metadata can be a very tedious and resource-consuming task, if it is not organised properly. One should avoid organising the capturing of metadata as a separate task. Instead the capturing of metadata should be integrated, as much as possible, with other tasks in statistics production, and this should be done in such a way that people see the meaning and importance of the task.

Ideally metadata should be generated as a side effect of other processes that have to be carried out anyhow. These processes will then become natural sources of metadata. The processes that generate metadata as by-products may be inside or outside the statistical system under consideration.

An important rule of thumb is that metadata should be captured as early as possible in the process of developing, implementing, and operating a production system. Another rule of thumb is that the same metadata should not be captured more than once.

Ideally metadata should be automatically transformed in parallel with the data that it describes. To the extent that this can be done, metadata production – e.g. the production of text labels and descriptions – will be automated in much the same way as the data production, the production of the figures, has been since long.

A major strategy for ensuring the overall efficiency and usefulness of the metadata-related work in a statistical organisation is to develop a global architecture for all metadata systems that the organisation may possibly develop during the foreseeable future. There are several purposes of such an exercise. One is to make an inventory of all metadata systems that need to be considered. This inventory should be developed in active cooperation between different stakeholders, so that after the exercise has been carried out, one can be reasonably sure that nothing important has been forgotten. Once developed, the inventory will be an excellent basis for putting priorities to different metadata-related tasks. Furthermore, by analysing and structuring the inventory, one may find that there are common components that can be used for different purposes, by different systems, provided that they are designed in a standardised way and are able to communicate via standardised interfaces. Such components may be given a higher priority from a global perspective than they should have got if looked upon from a more local point of view.

The development of "a grand plan" for the metadata infrastructure of a statistical organisation does not imply that this plan should necessarily be implemented in all its parts and in exactly the form originally specified. The plan should give visibility, clarity, and stability in the development efforts, but a lot of things will no doubt change during the implementation of the plan, which may take several years. Certain parts may never be implemented, other parts may be implemented in a different way than originally planned, and completely new components may appear as the result of new needs or new technical and methodological break-throughs.

A. Matching metadata needs and supply potentials in a statistical system

Every process in a statistical system has a need for different kinds of metadata. But every process has also a potential as provider of metadata. An important task when designing metadata systems is thus to match metadata needs and supply potentials in such a way that every process in a statistical system can get as many of its metadata needs as possible satisfied by natural metadata outputs from other processes in the system. This may require some adaptations in both metadata requirements and metadata outputs.

Table 1 indicates how the processes in a statistical system can be analysed from the following perspectives:

1. Which are the metadata needs of the process?
2. Are there any sources from which the process could get these metadata?
3. Which metadata could the process make available to other processes?

If this kind of analysis is carried out in some more detail for all major processes in a statistical system, one will obtain a good basis for rationalising the metadata processes of the system. Wherever a process could produce metadata that could satisfy a metadata need of another (or even the same) process, this should be done. In order to obtain "perfect matches", one may have to co-ordinate the specifications of the metadata concerned. This may sometimes require compromises, and both the metadata-producing and the metadata-using processes may have to be designed in slightly different ways from what was originally planned.

B. Designing sharable resources

This brings us to the issue of sharable metadata resources. We have just seen, how we may economise with metadata-related work by recircling metadata between processes within (and between) statistical systems. We may go further and create autonomous metadata resources that can be used by many systems and processes over a long period of time. Such sharable resources are most useful if the metadata they contain are designed in a standardised way, and if they are operated and maintained as separate tasks, with their own managements.

The sharable metadata resources of a statistical organisation would include different kinds of metadatabases (maybe with names containing words like "archive", "repository", "library", "catalogue", "dictionary", etc). The sharable metadatabases could contain

- production system documentations (with references/links to other types of metadatabases in this list)
- contents and quality descriptions of statistical end-products (typically sets of macrodata or microdata)
- definitions and descriptions of object characteristics and statistical characteristics (with references/links to sets of statistical data)
- definitions and descriptions of variables, classifications, and value sets
- definitions and descriptions of populations (with associated registers)
- questionnaires and other types of measurement instruments
- sets of coding rules (with references to classifications and thesauri)
- thesauri
- sets of editing rules
- knowledge about theories and methodologies
- knowledge about specific methods (with references to algorithms and tools)
- experiences, process data, and evaluation reports from this system and other systems
- information about standards
- reusable data models and process models

<p>Usage process</p> <p>Metadata needs of this process</p> <ul style="list-style-type: none"> • Which statistical data and services are available? Are they relevant for my task? • Contents and qualities of available statistical data (macrodata and microdata)? <ul style="list-style-type: none"> ◦ statistical microdata and their definitions ◦ relevance, accuracy, timeliness, availability, comparability, coherence <p>Potential sources of metadata needed</p> <ul style="list-style-type: none"> • Planning, design, and construction processes → e.g. definitions • Operation and monitoring processes → e.g. non-response rates (→ precision) <p>This process is a potential provider of</p> <ul style="list-style-type: none"> • Data about usage and requests (both satisfied and not satisfied) for data and services • Data about user satisfaction with products and services 	<p>Operation process</p> <p>Metadata needs of this process</p> <ul style="list-style-type: none"> • Information about processes: specifications and instructions for inputs, outputs, procedures • Information about methods, tools, and resources available • Feed-back about process performance in terms of errors, timeliness, resource consumption, etc <p>Potential sources of metadata needed</p> <ul style="list-style-type: none"> • Planning, design, and construction processes → instructions, specifications, descriptions • Operation and monitoring processes → process data <p>This process is a potential provider of</p> <ul style="list-style-type: none"> • Process data
<p>Planning and design process</p> <p>Metadata needs of this process</p> <ul style="list-style-type: none"> • Information about user needs and other stakeholder requirements • Information about methods, tools, available data sources, and other resources • Experiences (from this system and other systems) <p>Potential sources of metadata satisfying these needs</p> <ul style="list-style-type: none"> • Knowledge bases: physical and electronic libraries, websites, etc • Business intelligence processes: using search engines, intelligent agents as well as more traditional research methods <p>This process is a potential provider of</p> <ul style="list-style-type: none"> • Specifications of processes (inputs, procedures, outputs), descriptions of resources, instructions 	<p>Management and evaluation process</p> <p>Metadata needs of this process</p> <ul style="list-style-type: none"> • Information about user needs and other stakeholder requirements • Information about the performance of productions processes and usage processes • Knowledge about methods, tools, available data sources, and other resources • Experiences from comparable systems <p>Potential sources of metadata satisfying these needs</p> <ul style="list-style-type: none"> • Knowledge bases and business intelligence processes • Data generated by or requested from production and usage processes <p>This process is a potential provider of</p> <ul style="list-style-type: none"> • Experiences • Revised user needs and other stakeholder requirements
<p>Research and development process</p> <p>Metadata needs of this process</p> <ul style="list-style-type: none"> • Knowledge • Experiences <p>Potential sources of metadata satisfying these needs</p> <ul style="list-style-type: none"> • Knowledge bases and business intelligence processes • Experiences from production processes in many systems • Experiences from usage processes in many systems <p>This process is a potential provider of</p> <ul style="list-style-type: none"> • Knowledge, methods, and tools 	

Table 1. Where could statistical processes get their metadata from, and which metadata could they contribute to other processes? Matching metadata needs and supply potentials in a statistical system.

There are many more relationships between the sharable metadatabases than those indicated by the explicitly stated references and links in the list above.

By factoring out similar metadata resources from individual processes and production systems, and by generalising and standardising these resources, and making them sharable, one may rationalise the metadata management of a statistical organisation considerably.

Another important way of rationalising metadata management is by integrating it with data management wherever relevant. Actually a lot of the metadata associated with statistical data should primarily be stored as an integral part of the data sets. For example, individual observation records should contain quality data about the observations, such as data about missing values (why they are missing, etc).

Figure 8 gives an overview of a statistical system, where data and metadata flows are well co-ordinated, where metadata are recirculated, and where sharable metadata resources have been established.

C. Golden rules

Sundgren (2003b) suggest some "golden rules" for the development and maintenance of metadata systems. These rules are quoted, concluding this paper.

If you are a designer...

- Make metadata-related work an integrated part of the business processes of the organisation.
- Capture metadata at their natural sources, preferably as by-products of other processes.
- Never capture the same metadata twice.
- Avoid uncoordinated capturing of "similar" metadata – build value chains instead.
- Whenever a new metadata need occurs, try to satisfy it by using and transforming existing metadata, possibly enriched by some additional, non-redundant metadata input.
- Transform data and accompanying metadata in synchronised, parallel processes, fully automated whenever possible.
- Do not forget that metadata have to be updated and maintained, and that old versions may often have to be preserved.

If you are the project co-ordinator...

- Make sure that there are clearly identified "customers" for all metadata processes, and that all metadata capturing will create value for stakeholders.
- Form coalitions around metadata projects.
- Make sure that top management is committed. Most metadata projects are dependent on constructive co-operation from all parts of the organisation.
- Organise the metadata project in such a way that it brings about concrete and useful results at regular and frequent intervals.

If you are the top manager...

- Make sure that your organisation has a metadata strategy, including a global architecture and an implementation plan, and check how proposed metadata projects fit into the strategy.
- Either commit yourself to a metadata project – or don't let it happen.
- If a metadata project should go wrong – cancel it; don't throw good money after bad money.
- When a metadata project fails, learn from the mistakes, and do it better next time.
- Make sure that your organisation also learns from other statistical organisations.
- Make systematic use of metadata systems for capturing and organising tacit knowledge of individual persons in order to make it available to the organisation as a whole and to users.

