**❚❚❚❚ Statistics Finland**

# Multichannel publishing of statistics (electronic publications and database) - Finnish experience

## 1. The consept of multichannel publishing

### 1.1. Changing over from a printed statistical system to an electronic one

Official statistics constitute a statisti cal system comprised of a collection of esse n-tial social statistics of high quality that are produced regularly and sufficiently fr e-quently and are nationally representative.

General social statistics constitute the main product of Statistics Finland. At t he moment Statistics Finland produces some 200 sets of statistics on 26 topics. Some of the statistics are compiled monthly or quarterly, some annually or less fr e-quently. New statistical data are comprehensively released on the Internet. New data releases number about 650 per year.

Statistical information has traditionally been reported, released and disseminated in the form of printed publications, which go back a long time. Statistics Finland's oldest statistical publication series have been produced for  over 250 years. The evolvement and development of the Internet has put the traditional way of repor t-ing, publishing and disseminating statistical information in turmoil. The traditional form of publishing where tables and the information needed in their i nterpretation are bound into one volume in a publication series of official statistics is disappea r-ing. NSOs are changing over to electronic dissemination of statistics, both as far as tables (databases) and the text analyses (publications) needed in their  interpretation are concerned.

What kinds of characteristics should electronic dissemination of statistical inform a-tion and the entire statistical system posses when the transition is made from printed to electronic dissemination? Statistical information m ust be exhaustively available in the service, it must be easy to find, published data must always be a c-companied by the metadata needed in their interpretation and direct linking to them must be possible.

### 1.2. Common Structure of Statistical Information ( CoSSI)[1]

To implement good Internet-service the actual statistical information (tables, publ i-cations) and the metadata needed in its interpretation must be put into electronic format. However, this alone is does not create a good service. Automatic conve r-sion of information for different dissemination channels, its archiving, multiple language versions and attachment of sufficient metadata to it are all problems that must be solved in order to create a good service.

To solve these problems a Common Structure  of Statistical Information (CoSSI) was developed. The point of departure in the CoSSI was an infological analysis of statistical information. The conclusion from the analysis was that although in practice the definition of statistical information has vari ed according to a given

---

[1] Heikki Rouhuvirta, Harri Lehtinen. Common Structure of Statistical Information (CoSSI);
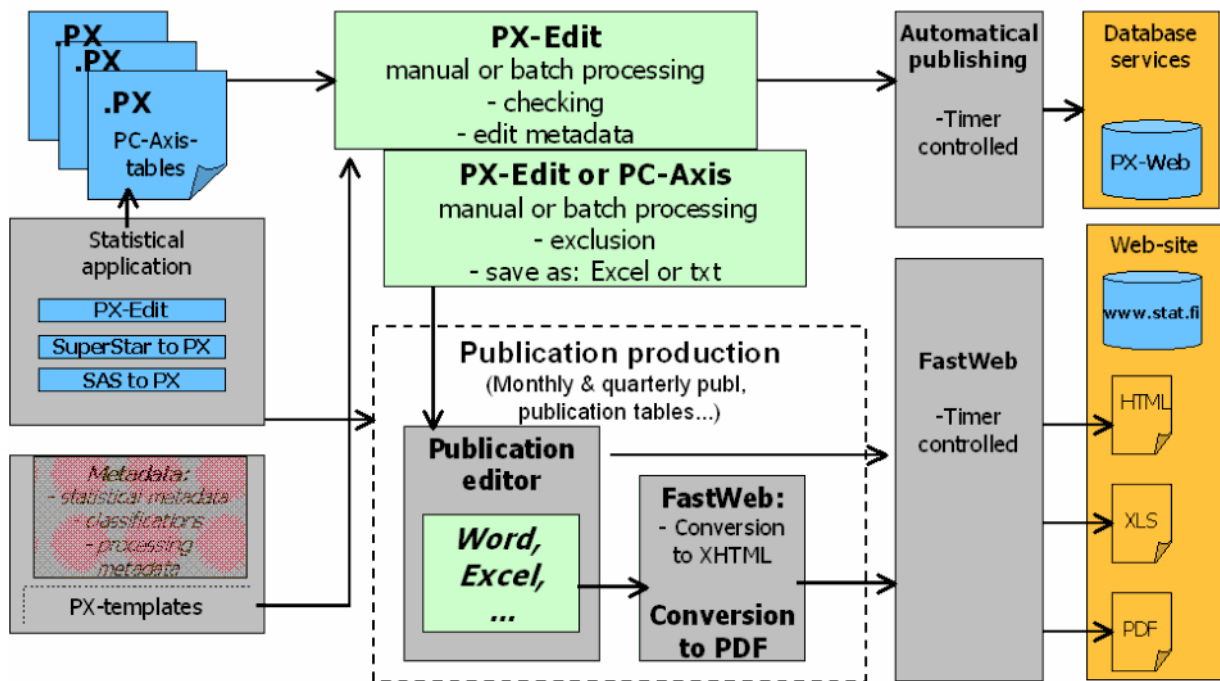http://www.stat.fi/cossi

situation and application, in reality statistical information has a certain simplifiable and acceptable universal structure. The CoSSI describes the general structure that is not dependent on the situation of the statistical informa tion presented in differing formats. CoSSI defines the structures of statistical data, metadata and publications.

The CoSSI model is a modular DTD system. It consists of Document Type Definitions (DTDs), it is based on standards (CALS, XDF, Dublin Core), it has a DTD for statistical matrixes and a DTD for statistical tables, it has DTDs for publications and documents. One XML file contains data, metadata and all language versions.

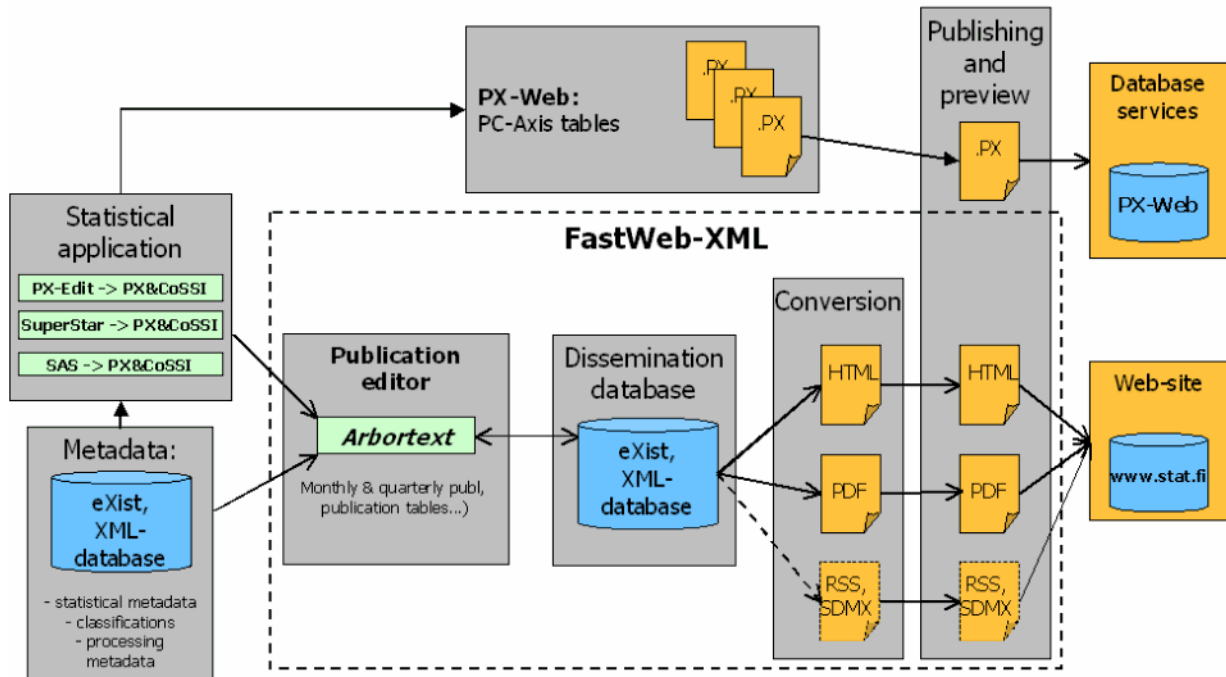## 1.3. XML publishing process compliant with the CoSSI model

In Statistics Finland's old publishing process, publishing from databases, updating of www pages and production of publications (printed, PDF and HTML) were all their own, separate production processes. Different tools were used in each process and overlapping work was done in them (Figure 1).

**Figure 1. Dissemination process – Office97[2]**



XML publishing process that complies with the CoSSI model will be fully impl e-mented in the dissemination of at Statistics Finland. In the new process, different formats of publications (printed, PDF, HTML) are generated automatically from one original XML file. Other forms of dissemination (email, RSS) are also pr o-duced automatically from the same original (Figure 2.).

---

[2] Still used in 40 statistics in May 2010; example "C onsumer Price Index";
http://www.stat.fi/til/khi/index_en.html

**Figure 2. XML based dissemination process - XML and PC-Axis[3]**



At the first phase, database tables (.PX) are produced with the same production a p-plications as the tables for the XML publishing process. SAS, SuperStar and PX - Edit produce both tables in XML format for the XML publishing process, and t a-bles in .PX format for database dissemination purposes (PX -Web). At this phase the preview functions and publishing for different publication formats (printed, PDF and HTML), and database tables (.PX) are combined within one transmission application and the files are saved into on e directory with a uniform structure. The same structure of topics and statistics becomes available in both the statistical www service and in the PX-Web database. The naming of XML, HTML, PDF and .PX files has been standardised at the different phases of the publishing process.

## 1.4. Archiving of an electronic statistical publication

Very little concern was given to the archiving of information published on the web when innovative services were being built in Internet's first decade of existence. The new situation has not been sufficiently examined by NSOs from the perspe c-tive of an entire statistical system in the long term either. How will the statistical i n-formation that is published today only in electronic format be available and acce s-sible in, say, ten years' time? What about in 30 years' time? If we now needed Co n-sumer Price Index data from 30 years ago, finding the data for 1980 would be po s-sible with the help of a publication series of official statistics.

In the production process that conforms with the CoSSI model a publication orig i-nal is saved in XML format in an eXist -XML database. This XML original of the publication contains the actual information content, the statistical metadata and the publication metadata in the desired languages in one XML file. This XML file is also archived into the XML database.

At the time of publication, intended dissemination formats (printed, PDF, HTML, etc.) are produced from this original XML file. The volume of metadata varies by publication format, for example, n umber of pages soon becomes a limiting factor in

---

[3] Used in 150 statistics in May 2010; example "Prices of Dwel lings";
http://www.stat.fi/til/ashi/index_en.html

printed publications. The most extensive metadata content can be attached to a pu b-lication in HTML format, whereby the desired metadata are printed out automat i-cally as HTML files in addition to the HTML fil es containing the actual statistical data.

The different formats (PDF, HTML) of publications published in the www service are published so that the URL addresses generated for different files at the time of publishing do not change later on. This makes usi ng the URL addresses elsewhere, such as diverse studies, articles or other web -based services, sensible and data user will always also be able to check the original document given in source references. Keeping published publications permanently accessible  to users is an essential ele-ment of a good information service and system of official statistics especially if data are no longer published in printed form. [4]

## 1.5. Roles of an electronic statistical publication and table database

Along with statistical databases, the volume of statistical information published in tabular format has grown and continues to grow enormously. However, tables pu b-lished in databases have a different role and different properties than tables pu b-lished in printed or electronic public ations. These differences deserve closer exa m-ining.

Tables published in databases are typically maintained as time series tables with continuously supplemented data contents. Database tables are also not archived in the same way as tables published in prin ted, HTML or PDF publications (or tables published in other file formats, such as Excel). For example, revisions of time series data, corrections of errors, or changes in classifications change the contents of a d a-tabase table retrospectively even in respe ct of earlier data, sometimes as long as years after its first publication.

Thus, in database tables, data on even older observations are always published a c-cording to the latest situation. Because of this, time series published in databases may also deviate from data relating to older points in time which have been pu b-lished as preliminary data in either a printed or an electronic publication or which have otherwise been revised or corrected subsequent to their first publication. For this reason, linking an electronic publication to the related database tables is not fully justifiable. The data only correspond with each other at the moment of pu b-lishing, for database tables may change later whereafter the data in them and in the electronic publication, and in the database no long match.

At Statistics Finland, the tables created from the PX -Web database are not linked to any other background database, but the PX -Web database builds up on the server from published PC-Axis files instead. In connection with the  introduction of the new XML publishing system each PC -Axis file published in the PX -Web server is automatically saved into a directory of statistics. Therefore, publications published on the home page of a set of statistics and the PC -Axis files saved into the same di-rectory can easily be permanently linked with each other. Then PDF or HTML fo r-mat publications published on the home pages of individual statistics and the r e-lated database tables (in PC-Axis format) become archived unchanged. This has two significant advantages: 1) it makes sense to also link the tables published in the database to the publication because they remain unchanged, 2) this way an archive is also created for the tables published in the PX -Web database.

---

[4] Statistical release archive; http://www.stat.fi/til/arkisto/index_en.html .

# 2. Multichannel publishing in practice in Statistics Finland

## 2.1. New practices for releasing statistics

Statistics Finland renewed its practices for releasing statistics in Finnish in 2004, and in Swedish and English in 2005[5]. According to the new publishing practices, a standard format release is produced for all statistics whenever new statistics are compiled and the data are made available to users. Before 2004, new statistics were published in press releases which numbered some 300 every year. The new publis h-ing practices have made the publishing of data more comprehensive. About 650 – 680 statistical releases have been produced annually, which shows the consi derable increase in the number of releases.

The new procedure was adopted for Swedish and English releases in June 2005. With the changeover, the number of Swedish language statistical releases grew from the earlier approximately 120 press releases to the same number as in Finnish, i.e. about 650–680 statistical releases, from the start of 2006. The number of En g-lish language statistical releases has increased from about 120 press releases ann u-ally to 350 statistical releases in 2010[6]. The objective is to increase the number of statistics released in English fu rther.

The modernisation of the publishing practices has taken into acco unt the objectives and requirements of the XML publishing process currently under construction from the perspectives of e.g. archiving and multichannel distribution. All statistical r e-leases are archived in the online service and their URL addresses will r emain un-changed.[7] The structure of statistical releases was carefully defined to meet the d e-mands of multichannel distribution. For example, the headings of statistical releases and the first paragraph are automatically delivered to the e -mail and RSS distribu-tions.[8]

Statistics Finland's XML-based publishing system became operational in the spring of 2007. The monthly publication of the Labour Force Survey, which served as the pilot in the XML project, has been created with an Arbortext XML editor since 22 May 2007 and published on Statistics Finland's online service in HTML and PDF format in both Finnish and English (with only the statistical release in Swedish), as well as in printed format in Finnish.[9] A single three-language XML original file can be used to automatically publish HTML pages in Finnish, English and Swedish (a total of about 90 HTML files), compile a PDF publication in Finnish and English, compile a PDF file in Finnish for printing and distribute the publication via e -mail and RSS in Finnish, English and Swedish.[10]

By May 2010, the publications of 150 of the good 200 Statistics Finland statistics were being created in both HTML and PDF format with the new XML process. So far the new publishing process has been used to produce over 1000 differe nt publi-cation titles in Finnish, English and Swedish (in both HTML and PDF format).

---

[5] Statistics Finland's new English web pages released (26 July 2005); http://tilastokeskus.fi/ajk/poimintoja/2005-07-26_webpages_en.html
[6] Release Calendar 2010; http://tilastokeskus.fi/ajk/julkistamiskalenteri/julkistamiskalenteri_aika2010_en.html
[7] Statistical release archive; http://tilastokeskus.fi/til/arkisto/index_en.html
[8] Latest statistical releases from Statistics Finland; http://tilastokeskus.fi/media/rss/2.0/tk_en.rss
[9] Content of the Labour Force Survey's home page is expanding; http://tilastokeskus.fi/til/tyti/tyti_2007-05-22_uut_001_en.html
[10] Labour Force Survey > 2010 > March; http://tilastokeskus.fi/til/tyti/2010/03/tyti_2010_03_2010-04-27_tie_001_en.html

## 2.2. Redesigning of publications

In the old publishing process (Figure 1), the data contained in a single set of stati s-tics was published 1) as electronic statistical rel ease in the online service (the co m-pulsory minimum; if required, annexed tables and figures and a longer article would also be published), 2) as a printed publication and 3) as database tables in the table database. All three publishing processes were sepa rate and employed different tools in the production of tables and the publishing of data. Data in different distr i-bution channels were not co -ordinated and they did not form a consistent whole.

The text section of the old printed Labour Force Survey was a two-column publica-tion in Finnish. The tables were bilingual (Finnish, English) and were optimised for A4-sized pages. The publication was available for a charge in both printed and PDF format. Alongside the printed publication, a statistical release and s ome annexed tables and figures produced in a different publication process and containing diffe r-ent data were released free of charge in the statistics online service on home page of the statistics.[11] The revised Labour Force Survey is a comprehensive publi cation composed of the text section, an extensive table and figure annex and a quality d e-scription.[12] The publication is available free of charge in both HTML and PDF fo r-mat, and a printed version of it can be o rdered for a charge.

The new publishing system (Figure 2) will combine the content of statistical r e-leases published in the online service (and the table and figure annexes supplemen t-ing them) and the content of earlier printed publications. This revision is not merely a question of routine copying of earlier data into new tools, but a complete rede s-igning of the content as a whole. The objective is to retain the publication form in the packaging of published data also in the f uture.

## 2.3. Basic publication from statistics

The role of publications being transferred over to the new publishing system is to release basic data on the official statistics. A basic publication is an entity that d e-scribes a single set of statistics at a certain reference time and it includes the stati s-tical release, a more exten sive text section, table and figure annexes and a standar d-ised quality description. If necessary, several basic publications from a set of stati s-tics can be released at a certain reference time, including a preliminary release and releases containing final data[13] or, for example, several publications focusing on different topics.

The objective is to transfer the basic publications from statistics comprehensively into the new XML-based publishing process. All 200 statistics in production will be converted to the new system by the end of 2010. In addition to technical changes (incl. the modernisation of table production so that they are produced in the XML format), the content of publications are redesigned as part of the changeover (texts, tables, general structure of the publication), the agency's publishing policy has been revised (charges for publications, whether or not a publication will be available in printed format, etc.) and the series division of official statistics has been r ethought.

## 2.4. Redesigning tables

The role of tables in a basic publication from a set of statistic is to describe the key points in a compact way. Even though PDF versions of publications are also created for multichannel distribution, the principle is that the electronic version is to be cre-

---

[11] Labour Force Survey 2007, March (statistical release in all languages, appendix tables only in Fi n-nish); http://tilastokeskus.fi/til/tyti/2007/03/index.html
[12] Labour Force Survey 2008, June; http://tilastokeskus.fi/til/tyti/2008/06/tyti_2008_06_2008 -07-22_tie_001_en.html
[13] Financial Accounts; http://tilastokeskus.fi/til/rtp/tie_en.html

ated first and foremost. The data content of a publication is optimised primarily for effortless use in an Internet browser. A publication's tables, therefore, must fun c-tion properly in HTML format, which requires that they be small and deal w ith only one subject. [14]

Although the rationale behind small tables is based on the obvious space and u s-ability restrictions, the use of small tables can be justified also on the grounds of i n-formation service. The message from usability experts is clear: T he users of Inter-net services want to access information on a website quickly and in an easily unde r-standable format. Large and complicated tables do not meet such information needs. The role of figures is also to depict the key phenomena in a compact way.

It has always been convenient to create figures to represent a single aspect, so there is no need for major changes on their part when publications are redesigned. In any case, the number of figures should be increased considerably, because their prope r-ties make them well-suited for online use. The lack of advanced graphics creation tools has been an obstacle to increasing their number. Such tools, however, are b e-ing sought at the moment.

**Figure 3. Three forms of a table: printed, A4, browser**

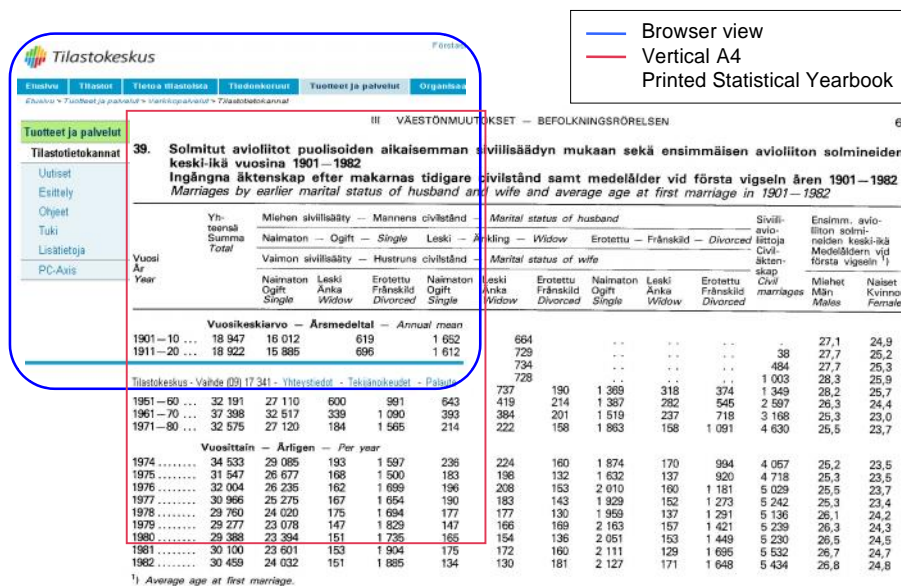| | Browser view |
| --- | --- |
| | Vertical A4 |
| | Printed Statistical Yearbook |

Figure: Jaakko Laakso

The role of database tables is significant when switching over from old -style printed publications to electronic publishing. Tables with extensive data contents must either be dismantled into small (browsable) tables dealing with one aspect or transferred to the database for publishing. In practice the majority of the tables in publications, which are larger than a single A4 -sized page, will be transferred to the database.

Therefore the basic publication from a set of statistics can be divided into electronic content published online, with two different formats and user interfaces (HTML and PDF), and tables published in a statistics database. The producer of a public a-tion faces a considerable challenge in compiling the publication in such a way that it functions properly in all of its formats and in different operating situations. This is a case of multichannel distribution in its full scale.

The availability of official statistics over time is greatly reduced if tables published in the database are not archived and s tored in a place where users can access them

---

[14] Telecommunications 2008; http://www.stat.fi/til/tvie/2008/ind ex_en.html

for a long period of time. For its table database, Statistics Finland uses the PX -Web database application, which does not have any links to any background databases; a PX-Web database is based on static PC -Axis files published on a server. Statistics Finland has started the archiving of different versions of database tables published in PC-Axis format, which will enable a publication representing a set of statistics at a certain time and the database tables rela ted to it to be permanently interlinked in such a way that the links remain unbroken when the publication and its database t ables are archived.[15] This will also enable the creation of permanent links between the compact publication tables and the large sour ce tables published in the database as well as their archived versions.

## 2.5. Interpreting data – metadata

Metadata – variables, concepts and definitions, classifications, quality descriptions, etc. – are essential to the correct interpretation and utilisat ion of data. In the past, the publication of metadata was restricted in printed publications by the amount of space available. The metadata published as part of a printed publication used to be easily available and remained so in archives as well. Various  more extensive meta-data publications that served as a user's handbook were also easy to archive. The space available for publishing metadata documents is no longer a restriction in the electronic format; all relevant metadata documents can now be published . Elec-tronic metadata documents should be equipped with the necessary IDs in the same way as the publications themselves in order to make the information easier to find.

A simple way of ensuring that the metadata required for the interpretation of data always accompanies the data itself is to publish both in the same package. This also applies to electronic publications. For this reason, among others, publications pr o-duced using Statistics Finland's XML publishing system shall always have a qua l-ity description attached to them. Matching each quality description version to the publication itself is easy when they have been combined in a single package. [16] The list of concepts and definitions is published through the concept database and corr e-sponds to the status of the latest publication; the archive is not available to the users of the information.[17] By attaching the concepts and definitions section to each pu b-lication it would be easy to ensure that the appropriate metadata always accomp a-nies the publication. On the other hand, classifications are often so extensive that it is not practical to attach them to a publication in their entirety. In such a case it is important that earlier classification versions are stored for access by users. [18] These should then naturally be also published in such a way that the entire classification can be stored at its original URL addresses and made available to users in the future as well.

## 3. Sources:

[1] Heikki Rouhuvirta, Markku Huttunen.  How NSOs can respond to changing user needs in the Internet era. Statistical Journal of the United Nations Economic Commission for Europe, Volume 20, Number 1 (2003), pp.  55-69.

[2] Heikki Rouhuvirta, Lehtinen Harri, Common Structure of Statistical Information (CoSSI) <http://www.stat.fi/cossi>.

---

[15] Population by gender and area 31.12.2007 and increase of population (archived PC -Axis file);
http://www.stat.fi/til/vaerak/2007/vaerak_2007_2008 -03-28_tau_101_en.px
[16] Migration, 2009; http://tilastokeskus.fi/til/muutl/2009/muutl_2009_2010 -04-22_en.pdf
[17] Migration, Concepts and Definitions;  http://tilastokeskus.fi/til/muutl/kas_en.html
[18] Industrial Classification - versions; http://tilastokeskus.fi/meta/luokitukset/toimiala/versio_e n.html