

## Detailed planning IT-report (Version 23 January 2018)

Id	Product	Deadline	Who*	Comments
8.3	<b>Report describing the IT-infrastructure used and the accompanying processes developed and skills needed to study or produce Big Data based official statistics.</b>	31-01-2018		
1.	<u>Metadata management (ontology)</u> It is important to have (high quality) metadata available for big data. This is essential for nearly all uses of Big Data. Ideally, an ontology is available in which the entities, the relations between entities and any domain rules are laid down.	10-01-2018	<i>Piet Daas (NL)</i> <i>Jacek Maślankowski (PL)</i> <i>Monica Scannapieco (IT)</i>	First draft is available
2.	<u>Big Data processing life cycle</u> Continuous improvement of Big Data processing requires capturing the entire process in a workflow, monitoring and improving it. This introduces the need to design and adapt the process and determine its dependence on external conditions.	10-01-2018	<i>Piet Daas (NL)</i> <i>Jacek Maślankowski (PL)</i> <i>Monica Scannapieco (IT)</i>	First draft is available
3.	<u>Format of Big Data processing</u> Processing large amounts of data in a reliable and efficient way introduces the need for a unified framework of languages and libraries.	20-12-2017	<i>Piet Daas (NL)</i> <i>Jacek Maślankowski (PL)</i>	First draft is available
4.	<u>Datahub</u> Sharing of multiple data sources is greatly facilitated when a single point of access, a so-called hub, is set up via which these sources are made available to others.	20-12-2017	<i>Piet Daas (NL)</i> <i>Jacek Maślankowski (PL)</i>	First draft is available
5.	<u>Data source integration</u> There is a need for an environment on which data sources, including Big Data, can be easily, accurately and rapidly integrated.	20-12-2017	<i>Piet Daas (NL)</i> <i>Jacek Maślankowski (PL)</i> <i>Sónia Quaresma (PT)</i>	First draft is available
6.	<u>Choosing the right infrastructure</u> A number of Big Data oriented infrastructures are available. Choosing the right one for the job at hand is key to assuring optimal use is made of the resources and time available.	20-12-2017	<i>Piet Daas (NL)</i> <i>Jacek Maślankowski (PL)</i> <i>Sónia Quaresma (PT)</i>	First draft is available
7.	<u>List of secure and tested API's</u> An application programming interface (API) is a set of subroutine definitions, protocols,	10-11-2017	<i>Jacek Maślankowski (PL)</i>	First draft is available

	and tools for building application software. It is important to know which API's are available for Big Data and which of them are secure, tested and allowed to be used.			
8.	<u>Shared libraries and documented standards</u> Sharing code, libraries and documentation stimulates the exchange of knowledge and experience between partners. Setting up a GitHub repository (or something similar) would enable this.	24-11-2017	<u>Jacek Maślankowski</u> (PL)	First draft is available
9.	<u>Data-lakes</u> Combining Big Data with other, more traditional, data sources is beneficial for statistics production. Making all data available at a single location, a so-called data-lake, is a way to enable this.	20-12-2017	<u>Piet Daas (NL)</u> <i>Jacek Maślankowski</i> (PL)	First draft is available
10.	<u>Training/skills/knowledge</u> Facilitating the exchange of skills and knowledge, for instance by offering trainings or assistance, will greatly stimulate the use of Big Data by many National Statistical Offices.	8-12-2017	<u>Piet Daas (NL)</u> <i>Jacek Maślankowski</i> (PL)	First draft is available
11.	<u>Speed of algorithms</u> Making use of fast and stable implementations of algorithms will greatly stimulate the application of Big Data as this speeds up the processing of large amounts of data tremendously.	13-12-2017	<u>Piet Daas (NL)</u> <i>Jacek Maślankowski</i> (PL)	First draft is available

\* alphabetical order and underlined main contributor, regular contributor and *italic "little" contributor*

## Detailed planning quality report

Id	Product	Deadline	Who*	Comments
<b>8.2</b>	<b>Report describing the quality aspects identified in studies focussing on the use of Big Data for official statistics</b>	<b>02-03-2018</b>		
1.	<u>Coverage</u> Information on the population included in a big data source is vital for reliable statistics. Important for this issue are the lack of information on the units included, their duplication and their selectivity.	02-03-2018	<i>Manca Golmajer and Crt Grahonja(SI), Jacek Maslankowski (PL) Tiziana Tuoto (IT)</i>	First draft available
2.	<u>Comparability over time</u> To produce comparable statistics over time, it is essential that the source remains accessible, relevant and its content remain usable.	02-03-2018	<u>Valentin Chavdarov (BG)</u>	First draft is available
3.	<u>Processing errors</u> During the processing of Big Data, errors may be introduced that negatively affect the quality of the data. Examples of this are the way outliers and missing values are treated.	02-03-2018	Manca Golmajer and Crt Grahonja(SI), <u>Magdalena Six (AT)</u>	First draft is available
4.	<u>Process chain control</u> In a Big Data process it is very likely that multiple partners are involved. To assure a stable and timely delivery of data of high quality, the entire process needs to be controlled.	02-03-2018	<u>Piet Daas (NL)</u>	First draft is available
5.	<u>Linkability</u> It is to be expected that Big Data needs to be linked or combined with to other data sources. During this process, errors may occur which affect the quality of the output.	02-03-2018	Manca Golmajer and Crt Grahonja(SI), <i>Jacek Maslankowski (PL) Tiziana Tuoto (IT)</i>	Draft of each chapter, consolidated within chapter-group: 26 <sup>th</sup> of January, Consolidated draft: 14 <sup>th</sup> of February Final: 2 <sup>nd</sup> of March
6.	<u>Measurement error</u> The values included in Big Data may not be all correctly measured; some may contain errors. This affects the outcomes produced, certainly when a systematic bias is introduced.	02-03-2018	<u>Valentin Chavdarov (BG)</u> Manca Golmajer and Crt Grahonja(SI),	First draft is available
7.	<u>Model errors and Precision</u> Big data based estimates are likely produced by models. The specifications of these models may be incorrect which	02-03-2018	<i>Manca Golmajer and Crt Grahonja(SI), <u>Magdalena Six (AT)</u></i>	Draft of each chapter, consolidated within chapter-group: 26 <sup>th</sup> of January, Consolidated draft:

	negatively affects the reliability of the estimates.		Tiziana Tuoto (IT)	14 <sup>th</sup> of February Final: 2 <sup>nd</sup> of March
--	--	--	--------------------	---

\* alphabetical order and underlined main contributor, regular contributor and *italic "little" contributor*

## Detailed planning methodology report

Id	Product	Deadline	Who*	Comments
8.4	Report describing the methodology (principles of finding and collecting Big data and assuring stable access, methodology of using Big Data as a single or major source of input, methodology of using Big Data as an additional data source in combination with others) of using Big data for official statistics and the most important questions for future studies	31-3-2018		
1.	<p><u>Assessing accuracy</u> How accurate are Big Data based the estimates. Both bias and variance need to be considered, but bias is expected to be more important.</p>	2-3-2018	<p>Tiziana Tuoto (IT)</p> <p>Manca Golmajer (SI)</p>	<p><b>Deadlines:</b></p> <p>- 24<sup>th</sup> of Feb 2018 – Draft Report</p> <p>- 2<sup>nd</sup> of March 2018 - Final report</p> <p>- 26<sup>th</sup> of Jan 2018 - Structure/template</p>
2.	<p><u>What should our final product look like?</u> Data driven statistics do not start with a predefined end product in mind. However, during this work it is important to start thinking about the product that can/will be delivered.</p>	2-3-2018	Valentin Chavdarov (BG)	<p><b>Deadlines:</b></p> <p>- 24<sup>th</sup> of Feb 2018 – Draft Report</p> <p>- 2<sup>nd</sup> of March 2018 - Final report</p> <p>- 26<sup>th</sup> of Jan 2018 - Structure/template</p>
3.	<p><u>Deal with spatial dimension</u> Many Big data sources have a spatial component, such as a geolocation,. It is essential to make use of this kind of information. This means that attention has to be paid to the location of objects and aggregating spatial data.</p>	2-3-2018	Piet Daas (NL)	<p><b>Deadlines:</b></p> <p>- 24<sup>th</sup> of Feb 2018 – Draft Report</p> <p>- 2<sup>nd</sup> of March 2018 - Final report</p> <p>- 26<sup>th</sup> of Jan 2018 - Structure/template</p>
4.	<p><u>Changes in data sources</u> Many Big Data sources are by-products of new technological developments. The content of these sources may therefore change rapidly. This will also be the case for data sources produced</p>	2-3-2018	<p>Valentin Chavdarov (BG)</p> <p>Piet Daas (NL)</p>	First draft available

	by private companies. It is important to get a grip on these to enable the production of reliable time series.			
5.	<p><u>Machine learning in official statistics</u></p> <p>Considering its rise in popularity, it is likely that in the future more and more (official) statistics will make use of machine learning based methods. It is vital to fully understand the implications of applying these kinds of methods in the production of official statistics. Estimation of variance, extracting features and knowing how to create good data sets for training purposes are examples of important considerations.</p>	2-3-2018	<p>Tiziana Tuoto (IT)</p> <p>Manca Golmajer (SI)</p>	
6.	<p><u>Data linkage</u></p> <p>To fully integrate Big Data in official statistics production it is essential that these data sources and/or their (intermediary) products can be combined with the data provided by other (more traditional) sources. Combining refers to the inclusion at either the unit or domain level here.</p>	2-3-2018	<p>Piet Daas (NL)</p> <p>Tiziana Tuoto (IT)</p>	<p><b><u>Deadlines:</u></b></p> <p>- <b>24<sup>th</sup> of Feb 2018</b> – Draft Report</p> <p>- <b>2<sup>nd</sup> of March 2018</b> - Final report</p> <p>- <b>26<sup>th</sup> of Jan 2018</b> - Structure/template</p>
7.	<p><u>Secure multi-party computation</u></p> <p>A lot of interesting data are produced by other organisations, such as private companies. Combining data from different organizations is challenging as many of them may not want the others to have complete access to their data; i.e. access at the individual record level. Being able to combine sources with a method that keeps the individual inputs of each partner private for the other parties involved is key to unlock the full potential of Big Data. National Statistical Offices are organizations that are ideally suited to function as a trusted party for this.</p>	2-3-2018	Piet Daas (NL)	<p><b><u>Deadlines:</u></b></p> <p>- <b>24<sup>th</sup> of Feb 2018</b> – Draft Report</p> <p>- <b>2<sup>nd</sup> of March 2018</b> - Final report</p> <p>- <b>26<sup>th</sup> of Jan 2018</b> - Structure/template</p>
8.	<p><u>Inference</u></p> <p>Methods that are able to reliably infer from Big Data and/or the combination of such data and other sources need to be developed. It is also essential to understand how this kind of inference is affected by the various sources of error that may occur.</p>	2-3-2018	<p>Piet Daas (NL)</p> <p>Tiziana Tuoto (IT)</p>	<p><b><u>Deadlines:</u></b></p> <p>- <b>24<sup>th</sup> of Feb 2018</b> – Draft Report</p> <p>- <b>2<sup>nd</sup> of March 2018</b> - Final report</p>

				- <b>26<sup>th</sup> of Jan 2018</b> - Structure/template
9.	<u>Sampling</u> Not all data are available or can be used. The ability to deal with subsets of Big Data and draw valid conclusions from it is important in unleashing its full potential.	2-3-2018	Valentin Chavdarov (BG)  Tiziana Tuoto (IT)	First draft is available
10.	<u>Data process architecture</u> Big data are often produced by one of the partners in the chain and subsequently transferred and/or used by others. For statistics production it is essential to have an overview of the whole chain and the individual steps performed by each partner as each step affects the other.	2-3-2018	Piet Daas (NL)	<b>Deadlines:</b>  - <b>24<sup>th</sup> of Feb 2018</b> – Draft Report - <b>2<sup>nd</sup> of March 2018</b> - Final report  - <b>26<sup>th</sup> of Jan 2018</b> - Structure/template
11.	<u>Unit identification problem</u> Big Data are produced by units of which (often) hardly any information is available in the source. This makes it challenging to identify the ‘real world’ units producing the data. For statistics it is essential to relate the units in Big Data with that of the (statistical) target population.	2-3-2018	Piet Daas (NL)  Tiziana Tuoto (IT)	<b>Deadlines:</b>  - <b>24<sup>th</sup> of Feb 2018</b> – Draft Report - <b>2<sup>nd</sup> of March 2018</b> - Final report  - <b>26<sup>th</sup> of Jan 2018</b> - Structure/template