# ESSnet Big Data II

**Grant Agreement Number: 847375**

https://webgate.ec.europa.eu/fpfis/mwikis/essnetbigdata
https://ec.europa.eu/eurostat/cros/content/essnetbigdata_en

## Workpackage L
## Preparing Smart Statistics

## Deliverable L4: Description of the findings regarding Task 4: Smart Traffic

**Final version, 1st November 2019**

### Prepared by:

Alessandra Sozzi (ONS, United Kingdom)
Henri Luomaranta (Statistics Finland)
Sverre Amdam and Johan Fosen (Statistics Norway)

Workpackage Leader:

Natalie Rosenski (Destatis, Germany)
Mail address:   natalie.rosenski@destatis.de
Telephone:      +49 611 754284

# Table of contents

# 1. Executive Summary

Task 4 examines data sources generated by traffic monitoring systems and vehicles that can be used within the context of official statistics. Subtask 4.1 outlines the work carried out by the Office for National Statistics (ONS - UK) and Statistics Finland (FI) on data collected from traffic monitoring systems (traffic loops data) and how this data is used to provide insights into economic activity. In subtask 4.2 Statistics Norway (NO) describes the overall landscape and availability of truck sensor data and how these data might be used in the production of freight transport statistics.

As outlined in subtask 4.1, both ONS and Statistics Finland have access to traffic loops data in the respective countries and in general we see a lot of similarities in the data and in the challenges to consider when handling the data: high frequency data, large size and high number of variables, challenge of missing data due to failures and dropouts of the sensors that compose the traffic monitoring system. One of the main advantages of traffic loops data is the availability of time series of traffic data often going back to times prior to the 2008 financial crisis. This allows us to test models in pre, during and post crisis times. ONS and Statistics Finland present different approaches for handling those challenges and targeting the economic activity.

In case study 1, ONS opts for an approach for creating a set of indicators that allow early identification of large economic changes. The aim is to provide insights into economic activity, rather than building official estimates of economic figures. We find that the road traffic counts give us some interesting insights into road transport in England. Although correlations are weak, the average traffic counts for the largest vehicles are consistent with at least some economic events (the financial crisis), and the trend broadly follows that of headline official economic statistics, although care must be taken with interpretation. These indicators are now released on a monthly basis since mid-April 2019 as research outputs.

In case study 2, Statistics Finland builds from the encouraging outcomes of the previous ESSnet Big Data and makes an extensive investigation using traffic loops data, relying on similar techniques used in the previous work. Using year-on-year growth rates of truck traffic as predictors and *a combination approach* of bundling together results from multiple machine learning algorithms, fairly accurate flash estimates of GDP can be produced, which reduce publication lag significantly. Previous ESSnet work has led to official publication of early estimates of Gross Domestic Product (GDP), both quarterly and monthly, as experimental statistics. The research outlined in this report demonstrates that traffic loops data can also be used for that purpose, with the additional benefit of data being available in almost real-time.

As a recommendation, indicators as well as early estimates of GDP are powerful tools to give an early picture of changes in the economy and may aid economic and monetary policy-makers and analysts in interpreting the economy. However, care must be taken when using them and it's important not to conflate them with our regular official estimates of GDP and trade. Those provide much more detailed numerical insights into how some economic activities have changed over a given period.

Moving on to subtask 4.2, data on road transport for official statistics is mainly collected through surveys nowadays. The existing survey design struggles with both quality and cost challenges, with high response burden, partial non-response and relatively high internal costs on data collection. This makes the search for alternative data sources particularly interesting in this field of official statistics.

A truck that rolls out of a manufacturer's fabric today contains about 2 000 sensors that measures everything from the technical performance of different engine parts, to the driving pattern of the driver. To get a better overview of different types of sensor data in trucks, we differentiate between three types of data; Fleet management system data (FMS) data, positioning data/GPS data and tachograph data. The three types represent different technologies and data, but also different data providers and possibilities of data access.

Following an exhaustive overview of each one of them, we looked at possible alternatives for utilizing sensor data in trucks for a "smarter freight transport survey". These possibilities for a modernised survey include: 1) using manufacturers' GPS-data and FMS-systems that allow the statistical office to access and collect the reported data in a standardized way, 2) integrating a web survey with a third part GPS-system or 3) integrating smart tachograph data in a modernized survey solution, if data can be remotely accessed.

## 2. Introduction

Task 4 examines data sources generated by traffic monitoring systems and vehicles that can be used within the context of official statistics.

With regards to 'connected cars', millions of SIM cards are already integrated in cars and trucks and are permanently sending and producing data, which could be used for subjects such as freight transport, population mobility or traffic accidents. Besides these smart devices integrated in vehicles, traffic monitoring systems have been in place in many countries for a long time, which also produce data that could be valuable for official statistics in several domains.

The task has been divided into two subtasks to address separately the two types of technologies. In particular, subtask 4.1 outlines the work carried out by the Office for National Statistics (ONS - UK) and Statistics Finland (FI) on data collected from traffic monitoring systems (traffic loops data) for the respective countries and how these have been used in production to provide insights into economic activity, whereas in subtask 4.2 Statistics Norway (NO) describes the overall landscape and availability of truck sensor data and how these data might be used to mitigate today's quality challenges and improve the production of freight transport statistics.

The report takes on each country's perspective but is also relevant from a European perspective as sensor technologies used in trucks and traffic monitoring systems are more or less the same throughout much of Europe.

In the remaining of this document we provide a detailed description of the two subtasks as follows:

**Subtask 4.1 – Use of traffic loops for economic estimates (FI, UK) –** ONS and Statistics Finland present two approaches of using data produced by road traffic monitoring systems to form indicators and early estimates of economic activity. The work done here builds upon the encouraging results obtained in Workpackage 6 "Early estimates" of the ESSnet on Big Data I (2016-2018)[1], where one of the most important contributions were the quite accurate nowcasts of monthly and quarterly gross domestic product (GDP) figures using firm data and some initial exploration of traffic loops. In particular, the work has been expanded in two ways:

1. **Faster Economic Indicators (UK) –** ONS constructed a number of indicators from traffic loops data obtained from Highways England[2], with an understanding of how they are related to economic activity, e.g. in the ability to forecast sudden shifts or long term trends.
2. **Nowcasting Finnish Real Economic Activity: a Machine Learning Approach (FI) –** Statistics Finland tested different models to nowcast real economic activity and GDP in near real time from traffic loops data obtained from the Finnish Transport Agency website[3], before the month or quarter of reference is over, in addition to the production of flash estimates.

---

1 See https://webgate.ec.europa.eu/fpfis/mwikis/essnetbigdata/index.php/WP6_Early_estimates1 for more details.
2 https://highwaysengland.co.uk/motorways/
3 https://aineistot.liikennevirasto.fi/lam/reports/LAM/

Both countries involved in subtask 4.1 have developed a set of processes and requirements in terms of methodology and IT systems for continued production of these indicators. These are currently used in production to support the regular publication of indicators and estimates.

**Subtask 4.2 – New data sources for producing statistics on road freight transport (NO) –** Statistics Norway investigates how widespread the use of smart devices attached to truck transport vehicles are in Norway. The main aim of this subtask is to obtain an overview of how different measurement instruments can contribute to freight transport statistics. This may include factory installed devices, monitoring devices from transport companies and prototypes which can be tailored for statistical purposes.

## 3. Subtask 4.1 - Use of traffic loops data in the context of official statistics

Traffic inductive loops sensors, or simply traffic loops, are a tried and tested technology for monitoring traffic. Sensors detect vehicles passing or arriving at a certain point, for instance approaching a traffic light or in motorway traffic[4]. In particular, the inductive loop is a square of wire embedded into or under the road. The loop utilizes the principle that a magnetic field introduced near an electrical conductor causes an electrical current to be induced. In the case of traffic monitoring, a large metal vehicle acts as the magnetic field and the inductive loop as the electrical conductor. A device at the roadside records the signals generated[5].

Both ONS and Statistics Finland have access to traffic loops data in the respective countries and in table 1 we outline the main differences and similarities between the two. This is followed by some considerations on the quality issues intrinsic to this type of data and how those were overcome by the two National Statistical Institutes (NSIs).

We proceed with focusing separately on the two different implementations of subtask 4.1, the *Faster indicators of UK economic* activity, adopted by ONS, and the *Machine Learning Approach for nowcasting Finnish Real Economic Activity*, adopted by Statistics Finland.

### 3.1 Traffic loops data: an overview from UK and FI

Traffic loops data is an easy to acquire data source, available to NSIs across several countries of the European Statistical System (ESS).

One of the main advantages of traffic loops data, when compared to other big data sources, is the availability of time series of traffic data often going back to times prior to the 2008 financial crisis. This allows us to test models in pre, during and post crisis times.

**Table 1: Traffic loops data comparison: England vs. Finland**

|  | ONS (UK) | Statistics Finland (FI) |
|---|---|---|
| **Provider** | Highways England[6] | The Finnish Transport Agency (Liikennevirasto)[7] |
| **Coverage** | England only | Finland only |
| **Reference period** | 2006 – up to date | 1995 - up to date |
| **Temporal aggregation** | 15 minutes | Continuous |

---

[4] https://en.wikipedia.org/wiki/Induction_loop
[5] https://www.windmill.co.uk/vehicle-sensing.html
[6] Access https://highwaysengland.co.uk/motorways/ for more details.
[7] Access https://vayla.fi/web/en for more details

| | | |
|---|---|---|
| **Frequency of update** | Available via Highways England API on the 20th of the following month. They contain 15-mins aggregated data from the previous month grouped by vehicle classes spanning the entire month. | Available at the 2nd day of each month. They contain data from the previous month grouped by vehicle classes spanning the entire month. Time stamp indicates the exact time of the measurement. |
| **Number of sensors** | ~13,000 sensors on Motorways, major trunk roads and junctions in England only | ~500 measurement points across Finland |
| **Types of sensors** | Three types:<br><br>- MIDAS: typically employed in rapid traffic flow management for example, smart motorways with variable speed limits.<br>- TAME and TMU: used for traffic monitoring, such as measuring bottlenecks around junctions and local road networks. | The TMS (traffic monitoring system) point consists of a data collecting unit and two induction loops on each traffic lane. The device registers vehicles passing the TMS point, recording data such as time, direction, lane, speed, vehicle length, time elapsed between vehicles and the vehicle class. |
| **Types of vehicles** | Vehicles are grouped by vehicle length:<br><br>• <5.2m: cars, motorcycles<br>• From 5.2m to 6.6m: panel vans, minibus<br>• From 6.5m to 11.66m: rigid lorries, buses<br>• >11.66m: larger rigid lorries and coaches, articulated lorries | Data are disaggregated by vehicle types:<br><br>• cars and delivery vans<br>• trucks<br>• buses<br>• semi-trailer trucks<br>• trucks with trailer<br>• cars and delivery vans with trailer<br>• cars and delivery vans with a long trailer or with a mobile home |
| **Data** | Traffic counts every 15 minutes by vehicle length and average speed over all vehicles. | Time, direction, lane, speed, vehicle length, time elapsed between vehicles and the vehicle class. |

As shown in table 1, the same type of data source presents sensible differences between the two countries. Additionally, the sensors that are part of the traffic monitoring systems are subject to failures and dropouts for long periods of time, which directly affect the number of sensors available at a certain time and the overall measurements as a result. Those can happen for various reasons, including (but not limited to):

• Roadworks
• Calibration test failures – if a sensor fails some calibration tests the data is not published. Could be taken offline for over a year.
• Hardware failures – again, could be taken offline for over a year.
• Software updates – they can be taken offline for updates that can last a few days.

- Power failures – some of these sensors are solar powered so may not have power all of the time
- Upload failures – data is lost if it cannot be uploaded.

**England vs. Finland data quality issues**

Along with these issues, the data used by ONS from 2006-2014 were published as aggregated to road sections. This is followed by a gap between December 2014 and April 2015 where a change in the data collection methodology took place. No data is available for this period. Since April 2015 traffic counts and average speeds were then collected for individual sensors.

To overcome such distortions the methodology used here is to construct three types of indicators for traffic counts, aggregated monthly for each vehicle size:

- Overall national (England) averages
- Overall average for sensors for all major ports
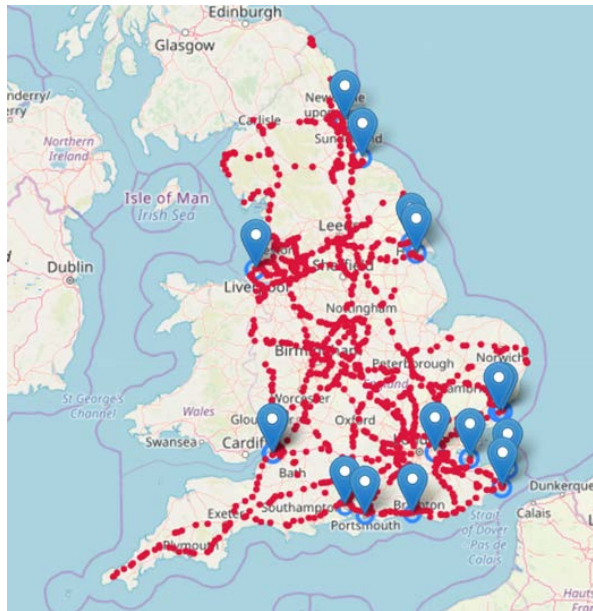- Averages for each individual port

Aggregates for ports were obtained by aggregating sensors within 10 km radius of 13 major English ports (a major port is defined as those that handle more than 1 million tonnes of cargo per year plus some historically important ports[8]). 10 km was chosen as this contains enough sensors to give a degree of robustness against sensor dropout. Some ports are further combined as their 10 km radius overlap. The published figures for all ports and the whole UK are finally seasonally adjusted (SA) using a standard JDemetra+[9] approach as this is robust to missing data. The steps are shown below in Figure 1b along with the locations of sensors and ports across England (Figure 1a).

---

[8] As defined in "Port Freight Statistics 2018: notes and definitions"
https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/735358/port-statistics-technical-note.pdf
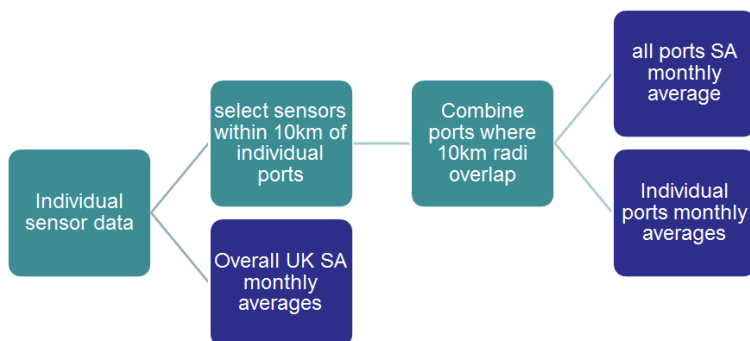[9] Software JDemetra+ https://ec.europa.eu/eurostat/cros/content/software-jdemetra_en

**Figure 1a: Location of road sensors**



Location of road sections with traffic sensors across England (red dots) and of the 13 major ports analysed in this work (blue pins and circles), in 2014. Note that the ports of Avonmouth and Royal Portbury near Bristol have been combined due to their proximity to each other. Felixstowe is combined with Parkeston for the same reason.

**Figure 1b**: **Data processing steps**



Compared to the ONS data, Statistics Finland used year-on-year growth rates to get rid of seasonality. Daily averages are used as the main unit of measurement aggregate at monthly level. Missing values e.g. due to malfunctioning sensor are imputed using generalized EM-algorithm by Josse & Husson (2016)[10].

---

[10] Julie Josse and François Husson. missmda: A package for handling missing values in multivariate data analysis. Journal of Statistical Software, Articles, 70(1):1–31, 2016. ISSN 1548-7660.

As the main caveat here, the study conducted by Statistics Finland used only Helsinki - Uusimaa region to compute estimates. The work can later be expanded to cover the entire country.

However it needs to be underlined that most of the economic activity is concentrated in this region, including the main ports. Adding more variables (and coverage) does not necessarily improve out-of-sample results. We were limited by feasibility considerations, as there is no Big Data environment in place, yet.

## 3.2 Subtask 4.1 - Case study 1: Faster indicators of UK economic activity

The Faster indicators of UK economic activity project[11], led by the Data Science Campus at ONS is an innovative response to the challenge of producing faster economic information. The goal of the project has been to identify several close-to-real-time big data and administrative data sources suitable for creating a set of indicators that allow early identification of large economic changes. The aim is to provide insight into economic activity, at a level of timeliness and granularity not currently possible with official economic statistics.

It is important to note that we are not attempting to forecast or predict GDP or other headline economic statistics here. Rather, the indicators should be considered as a means to provide an early picture of a range of activities that supplement official economic statistics and may aid economic and monetary policymakers and analysts in interpreting the fast changing economic situation in the UK.

Among the sources targeted as part of this project, the road transport sensors are one of them[12].

As highlighted in the previous chapter, we constructed average traffic counts and average speed indicators for the whole of England and for 13 major English ports, from traffic flow data from Highways England. Traffic activity around ports is of particular interest, as this may offer further insight into the understanding of the potential impacts of delays on trade and other economic activity.

We compare our monthly indicators of average traffic count and average speed with three official economic statistics:

- monthly gross value added (GVA), chained volume measure (CVM), seasonally adjusted (SA),
- monthly international trade in goods, imports (CVM, SA)
- monthly international trade in goods, exports (CVM, SA)

We make the comparison for the indicators for the all-England traffic flow and for the sum of all the English ports, for all vehicles and for vehicles disaggregated by length.

By overlapping the seasonally adjusted indicators with trade statistics and official estimates we find that the road traffic counts give us some interesting insights into road transport in England.

---

[11] See "Faster indicators of UK economic activity" https://datasciencecampus.ons.gov.uk/faster-indicators-of-uk-economic-activity/ for a more detailed overview of the project and latest data release.
[12] See "Faster indicators of UK economic activity: road traffic data for England" https://datasciencecampus.ons.gov.uk/projects/faster-indicators-of-uk-economic-activity-road-traffic-data-for-england/ for a more detailed description of the data, methods and results.

**Figure 2: Seasonally adjusted average all ports traffic counts, by vehicle length and international trade in goods (imports and exports, CVM, SA)**
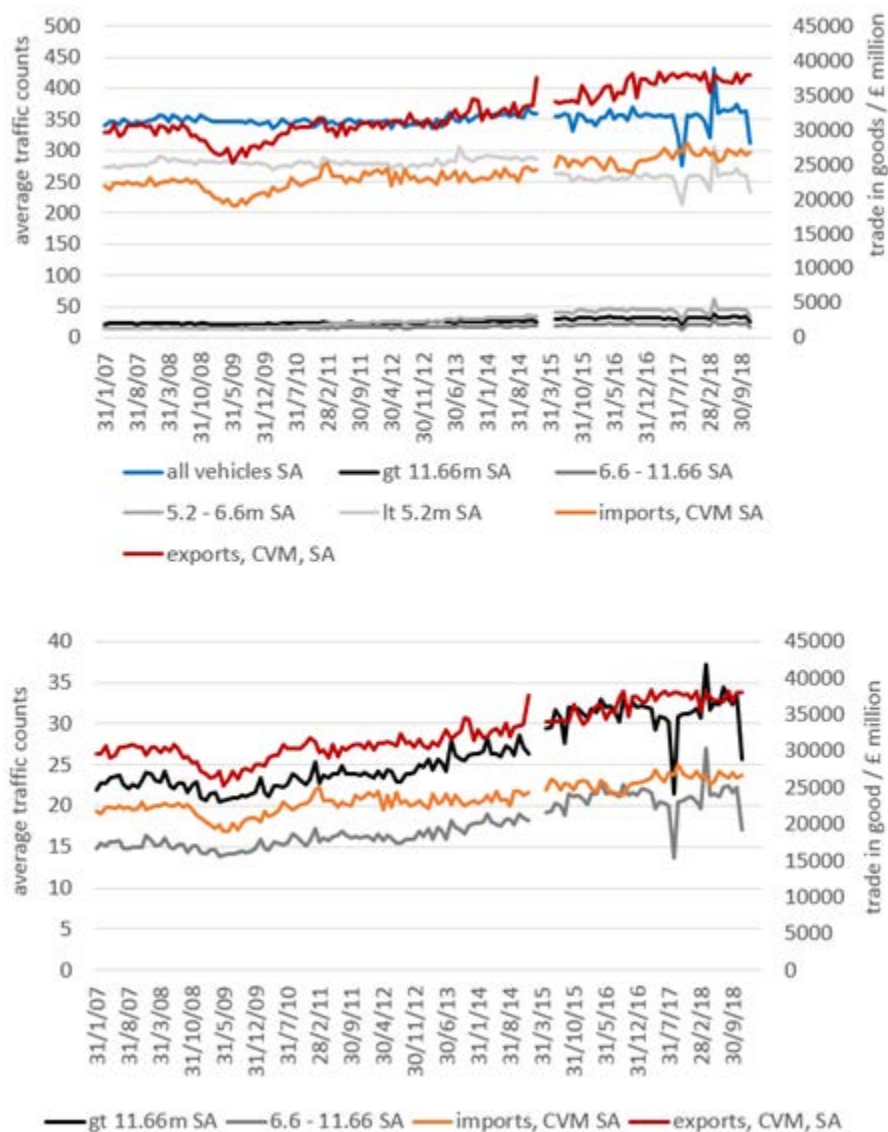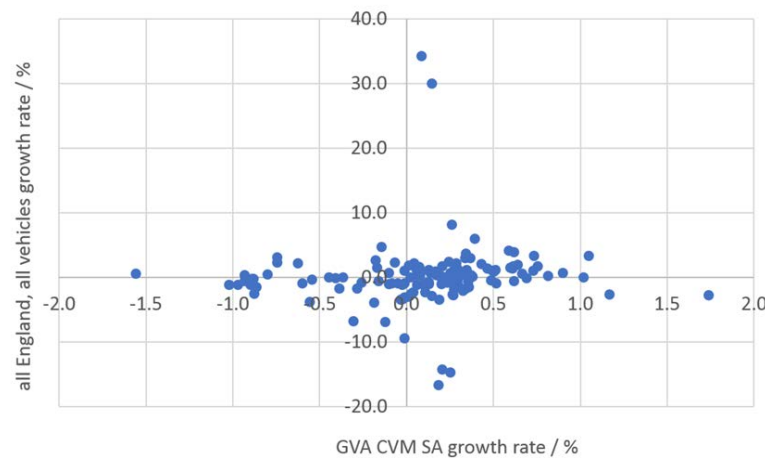




Figure 2 shows that in particular average traffic counts for the largest vehicles track imports and exports somewhat better than smaller vehicles. One might expect this to be the case, as international road freight is generally carried out using heavy goods vehicles rather than small cars. However, one might also expect the traffic counts to deviate from international trade estimates, as we make no correction or weighting for the value or volume of goods being transported, as we do not have this information. Movements in the data are consistent with at least some economic events (the financial crisis) and the trend broadly follows that of headline official economic statistics, although care must be taken with interpretation. Similar results have been obtained for ports traffic counts.

When we compare the change in traffic versus the monthly gross value added chained volume measure seasonally adjusted, we can spot outliers in traffic change which relate to spikes in data due to the effects of sensor dropout.  From figure 3, only a very weak correlation can be assessed – if any at all.

**Figure 3: all England average traffic counts for all vehicles (SA) and GVA (CVM, SA), growth rates, showing a large scatter in the relationship.**



Other measures, such as trade give similar results.

At this point the recommendation is that there is sufficient difference that these indicators should not be used to predict GDP on their own. Rather, they should be considered early warning indicators providing timely insight into real activities in the economy, and their ability to predict headline GDP should be carefully interpreted.

It may be that these indicators have the power to improve the performance of nowcasting or forecasting models, as components of these models, but we have not as yet tested this.

The quality of the underlying data represents so far one of the biggest obstacles for the correct and robust interpretation of the results. Although we tried to mitigate for sensor drop-outs by using the average, we haven't adjusted or modelled for any bias introduced by which sensors have dropped out. The drop outs can potentially have a large impact, and (under our current methodology) it's not clear what the impact is, which can make it difficult to interpret month-to-month changes.

The next steps will focus on improving the methodology and understanding of how these new indicators can be used as a faster indicator of economic activity, specifically, how traffic by different types of vehicles relate to local, national and international economic activity involving the transport of goods and people. However, this hasn't been pursued at present, due to other priorities.

*IT infrastructure*

The data is collected from the Highways England API[13] using a bash script to bulk download the data[14]. The data is downloaded and moved onto our in-house distributed system – the Data Access Platform (DAP), which is based on technology provided by the software company Cloudera. The platform is based on a powerful cluster of computers and provides many software tools to store the data in Hive tables then analyse the data in a pipeline built with Python and PySpark [15].

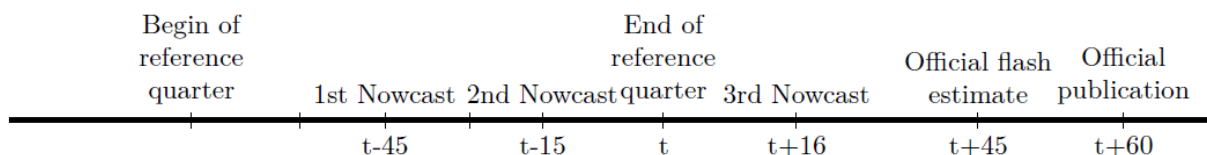## 3.3 Subtask 4.1 - Case study 2: Point estimates of Finland's GDP

Previous nowcasting exercises during the precedent round of the ESSnet Big Data[16], have considered firm level data in nowcasting framework to make early estimates of turnover indicators and GDP series.

The main lesson learned in the Finnish case, is that it is possible to create fairly accurate flash estimates and to reduce the publication lag significantly. This has been achieved by relying on firm level turnovers and *a combination approach* of bundling together results from multiple machine learning algorithms. This has led to an official publication of turnover indicators with a reduced lag and a publication of early estimates of GDP (quarterly and monthly) as experimental statistics. These releases are still being estimated by using firm-level data due to convenience, but underneath they rely on the improved techniques developed in WPL. The research presented below details how traffic loops data can be used to achieve similar results or produce even earlier estimates in almost real-time.

With WPL, a more extensive investigation is made using traffic loops data, improving upon the techniques used in the previous work. The purpose was to see whether one can do an even better job, given the timelines of traffic loops, in estimating point estimates of the Finnish economic output.

It is useful to summarize the timeline of events, corresponding to the table above:

**Figure 4: Timeline of events from the beginning of the reference quarter to the release of the official GDP estimates 60 days afterwards.**



We have computed the third month nowcast at t-15 days before the reference month has ended, simulating accurately the set of predictors available at that time[17]. In Finland, an official[18] flash estimate is provided 15 days before the official first quarterly GDP is released[19].

---

[13] http://webtris.highwaysengland.co.uk/api/swagger/ui/index
[14] Code to download the data is open sourced here: https://github.com/phil8192/webtri.sh
[15] Spark is a fast and general cluster computing system for Big Data. The Spark Python API (PySpark) exposes the Spark programming model to Python. See "Python Programming Guide" https://spark.apache.org/docs/0.9.1/python-programming-guide.html for more details.
[16] The reports of this work can be found in
https://webgate.ec.europa.eu/fpfis/mwikis/essnetbigdata/index.php/WP6_Documentation1

The following table describes the results that are obtained using traffic loops data and the combined approach. In appendix A we report the models tested and their respective references. The formulation of the exercise is done by extracting data from different traffic loops, and creating variables in such a way that each measurement point is an independent variable fed into the machine learning model.

**Table 2:**

|  | Nowcast second month | Nowcast third month | Nowcasts 16 days after | StatFi Flash |
|---|---|---|---|---|
| ME | 0.17 | 0.07 | -0.01 | -0.04 |
| MAE | 0.83 | 0.66 | 0.51 | 0.50 |
| RMSE | 0.99 | 0.85 | 0.66 | 0.64 |
| MaxE | 2.07 | 1.95 | 1.43 | 1.46 |

Table 2: ME (mean error), MAE (mean absolute error), RMSE (root mean squared error) and MaxE (maximum error) for the nowcast combination[20] approach, evaluated using the first version of quarterly GDP year-on-year growth that is officially released. Nowcast second month refers to the estimates of GDP computed during the second month of the reference quarter, nowcast third month are the estimates computed during the third month of the quarter and nowcasts 16 days after are computed after the end of the reference quarter. The results are from each quarter from the years 2014 to 2018. Statistics Finland's flash refers to the revision error of the current official first flash estimate, which is a useful benchmark. As a remark, the set of predictors in this table is entirely obtained from truck traffic and does not include any firm information.

The results indicate that it is possible to create real-time measurements using truck traffic as predictors (nowcast third month), with small sacrifice in accuracy. However, the increase in MAE is not dramatic, increasing only 0.16 %-points. Nowcasts 16 days after, provides already as good estimates as the current official flash estimate 1 month later.

Next steps include exploring more details on the aspects of capturing turning points in the economy using traffic loops. Even though the data we have shown so far is convincing in terms of providing accurate point estimates of the economic output in real time, more work is needed to understand how well we would capture dramatic fluctuations such as the Great Recession of 2009 in the Finnish economy.

The point estimates Statistics Finland regularly produces are a simple way of communicating the results, as opposed to, say, providing a confidence interval. Of course, the public may be taking the point estimates too literally, and be misled by the apparent exactness of the figures. This consideration needs to be carefully assessed when communicating the results and making statements. The Finnish practise is to provide the nowcasting errors every time a new release is made, so that the public can keep track on the reliability of the estimates in a transparent fashion.

---

[17] This means we are able to use exactly the data which would have been available to a practitioner.
[18] Official refers to the releases that are made with current methodology and are part of the regular "official" set of statistics produced by Statistics Finland.
[19] "Nowcasting Finnish Real Economic Activity: a Machine Learning Approach", the article summarizing all the details of these procedures, is found in: https://www.stat.fi/static/media/uploads/tup/nowcasting_empecon.pdf
[20] Combination is a simple unweighted average of the best 20 or so models, as detailed in the stand-alone article. We have tried other more sophisticated ways to do this, but find no improvements in the result.

*IT infrastructure*

All the Finnish exercises have been implemented using R software. In our opinion, R packages (such as Caret) are especially useful and R has proven flexible in both the testing and implementation phases.

Current statistical systems at Statistics Finland are implemented in SAS. Nonetheless, we have found it useful to carry out modeling in R, while data extraction and integration are done within the current SAS systems.

The statistical production system simply links to R routines that we have programmed and tested separately.

At Statistics Finland, only now we are starting a project to establish a proper machine learning and big data environment to facilitate machine learning solutions. Therefore, we cannot yet recommend an optimal way forward yet.

In the future, we plan to introduce a system where a cloud service is purchased from Microsoft (the Azure service).

## 3.4 Recommendations from subtask 4.1

Both ONS and Statistics Finland have decided to publish regular, monthly updates, using the most recent available data.

ONS has started to publish these new indicators based on traffic loops data on a monthly basis since mid-April 2019 as research outputs, while Statistics Finland has started publishing point estimates of GDP in January 2019 as experimental statistics based on firm level sales. Prior to that, since the end of 2018, nowcasts of turnover indicators are being released as official statistics. However, given the potential of traffic loops to provide even timelier estimates, the testing of this alternative data source should be continued.

Indicators as well as early estimates of GDP are powerful tools to give an early picture of changes in the economy and may aid economic and monetary policy-makers and analysts in interpreting the economy.

There has been a positive interest from the press[21] and follow-ups in social media, which more and more reinforce the idea that we are moving towards the right direction in responding to the challenge of producing outputs that meet the growing demand for more timely data.

While we believe this is already an important milestone, this is only the beginning of our journey in developing these products. We will continue to refine them to ensure that they can give as much useful information as the availability of data allows.

However, while these new products will hopefully give policymakers early warning of changes in the economy, it's important not to conflate them with our regular official estimates of GDP and trade. Those

---

[21] For the UK: some picked articles include Bloomberg, FT, The Times, City AM

provide much more detailed numerical insights into how some economic activities have changed over a given period.

# 4. Subtask 4.2 – New data sources for producing statistics on road freight transport

This section gives an overview of how data from truck sensors can be used for the future production of freight transport statistics. We proceed by describing the overall landscape and availability of truck sensor data taking a Norwegian perspective, and the potential utilization of these to produce future freight transport statistics. In appendix B we added some context for the reader about how data on freight transport for official statistics are currently collected, mainly through surveys, and what benefits sensor data can bring to the picture.

As mentioned before, despite the fact that subtask is mainly focusing on Norway, the report is also interesting from a European perspective as the road transport statistics are regulated in all EU and EEA countries with mostly the same information requirements. The data collection design is also similar and hence alike concerning quality challenges. Furthermore, the landscape of truck sensor data is likely to be similar across countries since technologies on sensor data in trucks is the same throughout much of Europe.

## 4.1 The concept of "Smart (truck) transport"

The concept of "smart transport" covers different topics where new technology and use of sensors is used for monitoring and improving efficiency in the field of transportation and logistics. Looking more specifically at "smart truck transport" it can be used to describe various concepts, ranging from using technology and sensor data to monitor and improve engine performance or to analye driving patterns, using fleet management information for better planning and efficiency of freight and passenger transport logistics to self-driving trucks. In all these concepts sensor data onboard trucks is the key, where large amounts of data are generated and stored, both from human-computer communication and non-human-communication or the Internet of Things (IoT).

A truck that rolls out of a manufacturer's fabric today contains about 2 000 sensors that measures everything from the technical performance of different engine parts, to the driving pattern of the driver. As an example, the total fleet of Volvo trucks in Europe (approx. 400 000 vehicles) generates around one petabyte of data each day. The use of sensors in trucks or truck fleets has been ongoing for many years, but earlier such data could only be accessed when the truck came into the garage. Today, onboard sensor systems transmit data in (almost) real time as the truck is on the road, by telecommunication to a cloud system or a remote server. Then it can be remotely accessed through dashboards, online management systems or shown as standardized reports provided by manufacturers.

## 4.2 Data sources - overview

To get a better overview of different types of sensor data in trucks, we can differentiate between:

- **Built-in sensor data**, or manufacturers' fleet management system (FMS) data, which is installed by the manufacturer before the vehicle is put into traffic. The manufacturer stores and processes the data.
- **Positioning data or GPS-data** that can both be pre-installed by manufacturers or third party devices installed in trucks by the transport enterprises, accompanied by an external FMS software .

- **Tachograph data** can be viewed as a source of its own, with a standardized setup across truck types, manufacturers and countries, containing standardized data on the driver, stop/start time, date etc.

### 4.2.1 Built-in sensor data from manufacturers

**Purpose of data**

A main purpose of the built-in sensor data is for the manufacturer and further authorized workshops, to monitor the performance of a truck and its different parts for technical, production and maintenance purposes. The manufacturers also have a substantial profit in processing performance data and selling them to truck owners/ transport enterprises, either as packages with dashboards, online management systems or standardized reports. The overall purpose here is FMS, where the transport enterprise management can monitor, evaluate and better plan the logistics and performance of the truck fleet.
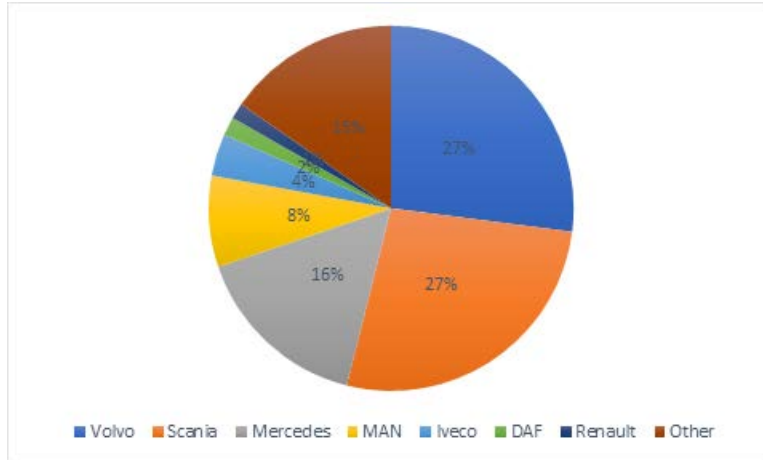
**Data providers**

All major truck manufacturers provide their own FMS to customers, and although there is a variety in how the solutions are built up, and to a certain degree what data can be disclosed in the solutions, much of the data from these systems are also standardized. Each manufacturer stores and processes the data. The seven largest truck manufacturers active on Norwegian roads all have their own FMS that are available for truck enterprises. Below is a list of all major manufacturers' FMS-systems with links to further information on each system:

- Volvo – Dynafleet
- Scania - Scania fleet managment
- Mercedes – Fleetboard
- MAN – MAN digital services (provided together with a third party FMS-provider called RIO)
- DAF – DAF connect
- Iveco – Iveconnect fleet
- Renault - Optifleet

To give an overview of the distribution of different truck manufacturers that operate in the transport sector, the figure below shows the market share of the seven largest truck manufacturers in Norway in 2018.

**Figure 5: Composition of the transport sector by truck manufacturers. Source: Statistics Norway**



## About the data – the "FMS standard"

- The data generated by built-in truck sensors are standardized across all major truck manufacturers through the so called international "FMS-standard" (http://www.fms-standard.com). The list of data that can be exchanged through the FMS-standard interface is comprehensive. However, the actual variables provided on the interface of each truck manufacturer depends on manufacturer, model and production year. All trucks produced after 2012 are provided with the FMS-standard. The main data that can be made available is:

- Vehicle speed (wheel based)
- Vehicle speed (from tachograph)
- Accelerator pedal position (0–100%)
- Clutch switch (on/off)
- Brake switch (on/off)
- Cruise control (on/off)
- PTO (Status/Mode)
- Engine coolant temperature
- Total fuel used (litres since lifetime)
- Fuel level (0–100%)
- Engine speed (RPM)
- Axle weight (kg)
- Total engine hours (h)
- FMS-Standard software version (supported modes)
- Vehicle identification number (ASCII)
- Tachograph information
- High-resolution vehicle distance
- Service distance

(For a complete list of data that can be provided through the FMS-standard see www.fms-standard.com)

Many of the above variables are not relevant for road freight transport statistics, but they are still mentioned above, since they might be useful as additional variables for other statistics in the future.

## Potential for official statistics

Most of the above listed data is not relevant for official statistics on freight transport based on today's statistical requirements in the freight transport survey, as these data are mainly monitoring driving behaviour and performance of different parts. But certain elements can be of interest. Axle weight is the most interesting as it gives crucial information on how much the truck is transporting, and separates journeys with cargo from journeys without cargo, which is an important distinction in today's survey.

This information could be useful as auxiliary information in estimation as well as in the sampling frame. As an example, the auxiliary information on distance contains, at least approximately, the same information as one of the variables of interest in the freight transport survey. Then, the auxiliary information can be used for quality evaluation of the variable of interest, although a limitation here is to what extent the auxiliary information covers the population of interest. Also it could be used as integrated pre-filled data in a modernized survey solution, combining new sensor data and survey data in a single data collection design. This will be elaborated in a later chapter.

Unfortunately, data on axle weight has proven to be unavailable on an initial analysis of FMS-data done in 2019 by the Norwegian economic transport institute (TØI), with a dataset provided from Volvo and Scania trucks. This should be further investigated to give a better understanding of the possibility of disclosure of axle weight data and for which manufacturers and models this can be done.

Other data that might be useful in the context of road transport statistics is the tachograph data, which gives data on engine start and stop. Tachograph data will be treated in a later chapter, but the FMS-data from trucks provides a link to tachograph data. Tachograph data in turn provides a link to data on the truck's register plate and thereby a link the national vehicle register.

## remote FMS-standard (rFMS)

In recent years some manufacturers have developed universal FMS, across types of brands and models. Thereby, they can offer the same data and reports in their FMS-solutions independent of truck brand or model. This is especially attractive for fleet owners who have a variety of different truck types and brands in their fleet. The technology used for this is the so called "remote FMS-standard", where a standardized rFMS-API is used to extract data from other manufacturers' FMS-sensors without the need to install any external hardware. All major truck manufacturers have agreed to make their sensor data accessible through the rFMS standard. Today, Volvo, Scania and Renault offer this technology in their FMS-solutions. Mercedes also have a universal solution in their FMS-system that is manufacturer independent and called Fleetguard, but their solution demands the installation of an external hardware device in all non-Mercedes trucks.

The existence of a universal system as above with an API-interface, means easier technical access to these data for future official statistics.
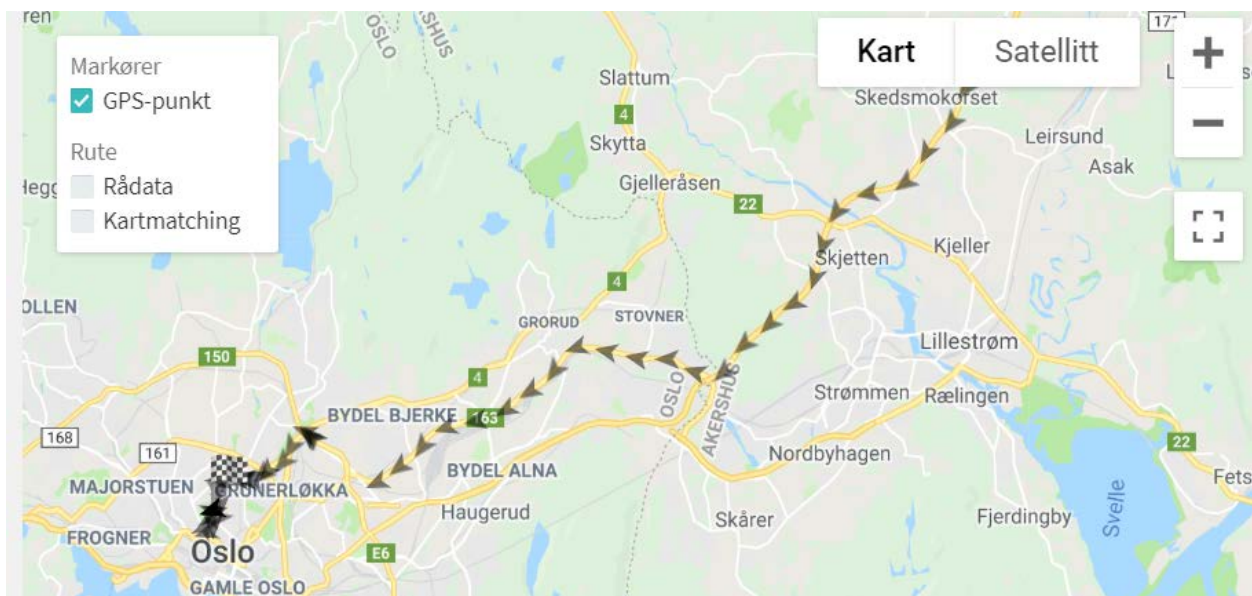
### 4.2.2 Position data

**Purpose of data and data providers**

Position data or geodata, using GPS-technology is not included in the FMS-standard, and thereby not as standardized as the built-in sensors in the above mentioned FMS-data. But since position data is a key element for FMS, it is standard to have GPS-technology preinstalled in the trucks by manufacturers. For newer trucks, all major manufacturers provide positioning as a feature in their FMS-dashboard systems. The built-in GPS-data can also be made accessible for other manufacturers' FMS-systems (and others interested), through the earlier mentioned rFMS API. This is an interesting gateway to GPS-data. It also makes the data somewhat more standardized across the manufacturers. Today Volvo, Scania and Renault have opened for this. Nevertheless, every manufacturer has their own GPS tracker technology installed with a potential variety on updating frequency rate on vehicle position (data points), and methods of processing the data into map-solutions and data visualization. Moreover, there are numerous third party FMS-providers that offer positioning data in their own FMS-systems. This is done through separately installed GPS-hardware in each truck. There are different technologies and patents used, such as stand alone GPS-sensors, and GPS-sensors connected to OBD2 ports (On Board Diagnostics). Some of the latter also have access to the truck sensors through the rFMS API. There are several providers on the Norwegian market, both domestically developed solutions and international solutions.

**About the data**

Positioning data or geodata are generated from a GPS-tracker containing latitude, longitude, altitude, driving direction and timestamps with a relatively high updating frequency depending on the movement of a vehicle. Data is usually generated if the truck is moving and halts when it stops. To make sense of the GPS-data they are processed in a geographical map-layer, such as google maps or similar. In FMS systems, GPS-data is usually displayed as dots or lines following the route of the truck.

**Figure 6: GPS data displayed as data points on a google map - from third party FMS-system**

Most FMS-solutions also register start address and stop address for each registered journey to differentiate the journeys undertaken. Below is an example from a third party FMS-system tested in a transport vehicle by Statistics Norway in 2018. Each journey is separated in the FMS-solution, with start and stop address, distance travelled, time used and mapping of the route, based on the coordinates of the GPS-data points registered on the particular trip.

**Figure 7: Illustration of how journey data are registered in a third party FMS-solution**



## Potential for official statistics

Positioning data provides information and insight on truck movement that might be very valuable for statistics on freight transport. Much of today's survey questionnaire is about positioning all journeys of a sampled truck in the reporting period with start and stop address information. Using GPS-data as a supplement to or a replacement for manually reported position data, and maybe also incorporating it in the survey questionnaire, could be a way of increasing data quality and lower response burden as well as survey costs.

By using GPS-data as auxiliary information in estimation as well as in the sampling frame, one could use this data as a supplement data source to the freight transport survey, potentially enriching some aspects of the data quality. More ambitious is using GPS-data as a replacement data source for position reporting in the survey. A way of doing this is prefilling of journeys in the web survey solution, based on the GPS-data, where respondents can validate the prefilled information in the web questionnaire and fill out auxiliary information on the goods transported of each journey registered. This can potentially lower response burden, lower partial non-response and also make the data more precise as the GPS-sensors will capture all journeys and structure them in a standardized way.

A challenge when dealing with the journey registrations in the trucks GPS-data is incorrect registration of journey stops. Initial testing of an FMS-system done by Statistics Norway and data analysis done by the Norwegian transport economic institute, found that GPS-sensors tend to interpret shorter stops, for example stops on a red light or traffic stops in congestion, as the journey end point. Journeys also tend to be registered when the truck is driving a few meters, for example when loading or unloading. This is a quite frequent error that creates a form of noise in the data, which we must deal with when the data is used as a part of a survey.

### 4.2.3 Tachograph data

**Purpose of data and content**

The tachograph is a device that records the driving time, breaks, rest periods as well as periods of other work undertaken by a driver. The purpose of the tachograph is to enforce the rules on driving times and rest periods and monitor the driving times of professional drivers in order to prevent fatigue and guarantee fair competition and road safety. From 2006 all vehicles above 3,5 tons are obliged to have a digital tachograph installed in all EU and EEA countries (European commission 2019).

The tachograph and the data recorded are standardized throughout the EU/EEA. Truck drivers/transport enterprises are obliged to read out tachograph data from the devices frequently for control and archive purposes. To read out tachograph data there are numerous software solutions that process the data into standardized reports. FMS-systems, both from manufacturers and third party, often have solutions where they extract tachograph data remotely and process it as a feature in their FMS-systems. Today's digital tachograph register the following data:

- Date Vehicle registration number
- Vehicle speed (for previous 24 hours of driving)
- Single or co-driver
- Number of times a driver card is inserted each day
- Distance travelled by the driver – it reads the vehicle odometer when the card is inserted and removed
- Driver activity – driving, rest, breaks, other work, periods of availability
- Date and time of activity change
- Events (for example driving without a driver card, overspeeding, fraud attempts etc) and faults
- Enforcement checks

**Smart tachographs**

By June 2019, a new regulation from the European council came into force, making a new generation of *smart tachographs* mandatory in new licenced trucks in all of Europe, and thereby setting a new standard for tachographs and the data these devices generate.

The smart tachographs have integrated positioning technology (GNSS and GPS) in each device, with standardized data recordings. That makes it easy for external systems to access the tachograph data for further processing (Intelligent transport systems – ITS). Not much information on the data content and technology can be found, but all specifications are detailed by the European commission in the current regulation (European commission 2019).

**Potential for official statistics**

On one hand, the earlier generation of tachograph data generates little data of direct interest based on the information needs for official statistics on freight transport. The introduction of the smart tachograph, on the other hand, might be very valuable for future road transport statistics as it records positioning data, and makes data access much easier. Data will come from one single system technology and structured in the same way across manufacturers and countries. It will also reach a level of full

representativity for new trucks on the roads as the tachograph data will be mandatory in every truck independent of manufacturer, type of truck, goods transported and geography. In addition, tachographs and its data is not market driven, but governmentally regulated and upheld by national transport authorities. This can make it easier for other governmental institutions to access and utilize these data, for example for the National Statistical Institutes (NSIs) for making official statistics.

An obvious downside with the smart tachograph data is that it is only mandatory in newly registered trucks from June 2019 and forward. This means it will take many years before this technology is widespread through a large enough share of the total truck fleet in Norway (and rest of Europe). Nevertheless, further exploration of this data source for future statistics on road transport should be undertaken and perhaps an early involvement of NSIs and EU-institutions towards the institutions responsible for tachograph legislation and content could prove valuable for the future.
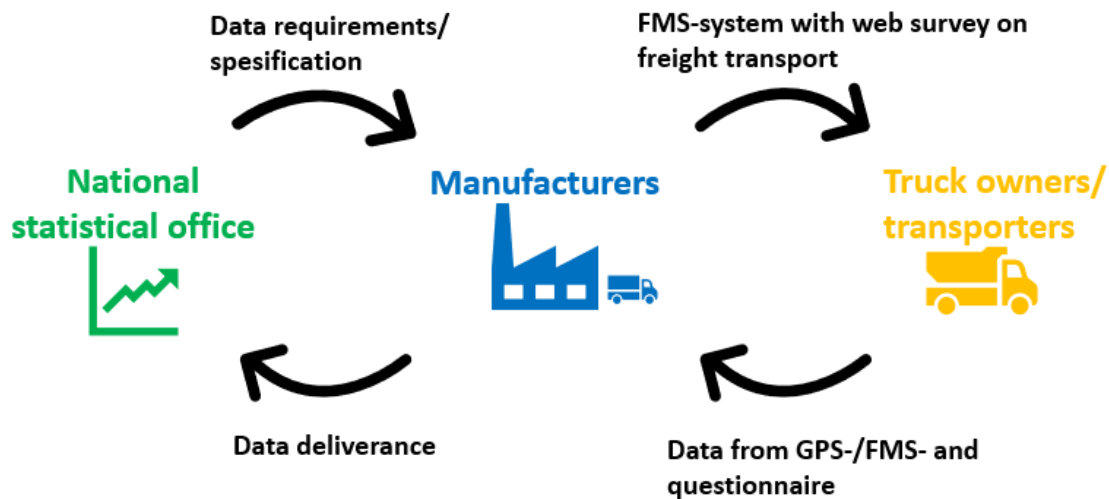
## 4.3 Recommendations from subtask 4.2

Below, we will look at possible alternatives for utilizing sensor data in trucks for a "smarter freight transport survey".

### Modernized survey using manufacturers GPS-data and FMS-systems

If one is to use GPS-data as a part of a modernized survey design, one possible source is the manufacturers built-in GPS sensor data, which will require a form of collaboration with each manufacturer. A possible solution is to incorporate a web survey module in the manufacturers FMS-systems/dashboards as a feature for truck owners/transport enterprises, where they can report to the statistical office if sampled in the survey. Here, one could also incorporate axle weight data from the FMS-sensors if available. With both position data and axle weight data on each journey, only information on the goods transported would be left for the truck owner/enterprise to register manually. By using the rFMS-API, the statistical office could access and collect the reported data in a standardized way.

The upside to this solution is that there is a business case for all parts. The truck owner/enterprise would benefit such a solution, as the response burden in the freight transport survey would decrease when the reporting would be far more user-friendly. The manufacturers would benefit by offering a new feature to their FMS-solution, making it more attractive for the truck owners to subscribe to their FMS-packages. And for the NSI, the benefit would be more data, better data and a more modernized and effective data collection design.

**Figure 8: Information model - modernized survey using manufacturers GPS-data and FMS-systems**



On the other hand, there are several hurdles and challenges to such a system. For the manufacturers it could prove both difficult and costly to establish and maintain such a solution. The NSI would lose control of the data collection to the truck manufacturers FMS-systems, where changes, updates and adjustments in the survey could be both costly and difficult to get done. As each manufacturer has their own FMS-system, it might demand the development of a separate survey solution for each manufacturer, creating methodological difficulties as well as making the development process costly and time consuming.

Despite all possible challenges, an attempt to go into dialogue with one of the largest manufacturers and possibly develop a demo should be investigated further.

**Modernized survey using third part GPS-data and FMS-systems**

An alternative to using manufacturers GPS-sensor data is integrating a web survey with a third part GPS-system, offering sampled enterprises/truck owners to install an external GPS-device in the sampled trucks for the reporting period. Numerous third-party FMS-systems exist on the market, offering a manufacturer-independent FMS-system, with their own GPS-hardware and a FMS-software system (web based and/or installed software). Statistics Norway have done some early exploration on such a system, testing functionality and features, with interesting and promising findings, but a lot more testing, data analysis and development should be done.

An advantage by using third party GPS-sensors in a survey context is the ability to create a web survey design and data collection in one system and customize it according to the NSIs standards and information requirements. The GPS-sensor's data recording could be tailored for statistical purposes, and also the web survey solution. A disadvantage is that NSIs have to rely on external hardware installation in each sampled truck. This could be costly, as well as introducing practical issues with shipping and installation, which would represent a different form of response burden.

**Modernized survey using smart tachograph data**

A third way to utilize position data in trucks in a modernized survey is to integrate smart tachograph data. If data can be remotely accessed from the tachographs, preferably with continuous transfer of data, one could imagine a web survey solution where position data from tachographs were directly transferred during or after journeys. Respondents could use an easy accessible portal (app or web portal) to validate journey data and manually add auxiliary information about the goods transported.

This solution will be of more relevance in years to come as smart tachographs become more common in the truck population, but early testing and exploration would give valuable insight for the development of such a design.

# 5. Conclusions

This report examines data sources generated by traffic monitoring systems and vehicles that can be used within the context of official statistics. In particular, subtask 4.1 outlines the work carried out by ONS and Statistics Finland on data collected from traffic monitoring systems (traffic loops data), whereas in subtask 4.2 Statistics Norway describes the overall landscape and availability of truck sensor data and how these data might be used in the production of freight transport statistics.

In subtask 4.1 we initially explore similarities in the data collected by the two NSIs and the challenges to consider when handling them: high frequency, large size, high number of variables and the presence of missing data due to failures and dropouts of the sensors that compose the traffic monitoring system. One of the main advantages of traffic loops data is the availability of time series of traffic data often going back to times prior to the 2008 financial crisis. This allows us to test models in pre, during and post crisis times. ONS and Statistics Finland presents here two different approaches to produce early indicators of economic activity: ONS opts for creating a set of indicators that allow early identification of large economic changes whereas Statistics Finland derives truck traffic growth rates to produce early estimates of headline economic figures.

We find that the road traffic counts give us some interesting insights into road transport in England. Although correlations are weak, the average traffic counts for the largest vehicles are consistent with at least some economic events (the financial crisis), and the trend broadly follows that of headline official economic statistics, although care must be taken with interpretation. At this point the recommendation is that there is sufficient difference that these indicators should not be used to predict GDP on their own. Rather, they should be considered early warning indicators providing timely insight into real activities in the economy. Their ability to predict headline GDP should be carefully interpreted. The quality of the underlying data, affected by sensor outages, represents so far one of the biggest obstacles for the correct and robust interpretation of the results.

Statistics Finland, on the other hand, builds from the encouraging outcomes of the previous ESSnet Big Data and makes an extensive investigation using traffic loops data, relying on similar techniques used in the previous work. Using year-on-year growth rates of truck traffic as predictors and *a combination approach* of bundling together results from multiple machine learning algorithms, fairly accurate flash estimates of GDP can be produced, which reduce the publication lag significantly. This has led to an official publication of early estimates of GDP (quarterly and monthly) as experimental statistics.

In subtask 4.2 the report describes the overall landscape and availability of truck sensor data, how these data might be used to mitigate today's quality challenges and improve the production of freight transport statistics. Today, data on road transport for official statistics is mainly collected through surveys. The existing survey design struggles with both quality and cost challenges, with high response burden, partial non-response and relatively high internal costs on data collection. This makes the search for alternative data sources particularly interesting in this field of official statistics.

In this report, we differentiate between three types of data: FMS data, positioning data/GPS data, and tachograph data. The three types represent different technologies and data, but also different data providers and possibilities of data access. Positioning data generated from built-in or third part GPS-sensors is the most promising data type with potential of replacing today's manual reporting of truck movement and improving data quality through a smarter survey solution, and also in

sampling/estimation models. These data are either provided by built-in manufacturer solutions or third party solutions with installed hardware. Furthermore, a new generation of tachographs, made mandatory in all European trucks from June 2019, generates and process positioning data in a standardized way with a large potential for future statistics. Also sensors, providing data on the trucks axle weight from manufacturer's FMS, might be of relevance to future freight transport statistics containing important information to the statistics, but with a more uncertain availability. We give some further reflections on how the new data sources might be incorporated in a new data collection design.

# 6. Appendix

## A. List of models for subtask 4.1 - Case study 2 (FI)

| Model/Technique | Name in caret | Reference |
|---|---|---|
| Factor models/principal components regression | pcr | Stock and Watson (2002a) |
| Independent component regression | icr | Hyvärinen. and Oja (2000) |
| Ridge regression | glmnet | Hastie et al. (2009), Chapter 3.4.1 |
| Lasso | glmnet | Hastie et al. (2009), Chapter 3.4.2 |
| Elastic-net | glmnet | Zou and Hastie (2005) |
| Least angle regression | lars | Hastie et al. (2009), Chapter 3.4.4 |
| Bayes Generalized Linear Model | bayesglm | Gelman et al. (2008) |
| Gaussian process | gausprLinear, gausprPoly, gausprRadial | Williams and Barber (1998) |
| Partial least squares | kernelpls, pls, simpls | Hastie et al. (2009), Chapter 3.5.2 |
| Bagged MARS | bagEarthGCV | Hastie et al. (2009), Chapter 9.4 |
| Regression Trees | ctree | Hastie et al. (2009), Chapter 9.2.2 |
| Boosting | BstLm, gbm, xgbTree | Hastie et al. (2009), Chapter 10 |
| Random forests | parRF, ranger, RRFglobal | Hastie et al. (2009), Chapter 15 |
| Nearest-neighbors | knn, kknn | Hastie et al. (2009), Chapter 13 |
| Neural network | pcaNNet | Hastie et al. (2009), Chapter 11 |
| Support vector machine | svmLinear, svmPoly, svmRadial svmRadialCost, svmSigma | Hastie et al. (2009), Chapter 12 |
| Penalized regression | penalized, rqnc | Hastie et al. (2009), Chapter 16 |

## B. Overview of statistics on transport of goods (NO)

All EU-member countries and EEA-countries are committed, through a union directive (1998), to deliver "comparable, reliable, harmonized, regular and comprehensive statistical data on scale and development of the carriage of goods by road" on a yearly basis (Eurostat 2017). This regulation has been put in place to provide both the EU and national governments necessary information for monitoring, controlling and evaluating road transport both inland and throughout the union.

Today, data on freight transport for official statistics is mainly collected through surveys. The main pillar for producing statistics on road transport is the Eurostat-regulated Survey on freight transport.[22] This is a comprehensive survey with the goal to provide representative statistics on all types of road transport carried out by freight trucks with load capacity above 3,5 tons. In Norway this survey provides the main statistical product for monitoring transport of goods on Norwegian roads.

Today's survey design struggles with both quality and cost challenges. A comprehensive and complex web survey questionnaire makes the response burden high, affecting partial non-response and the quality of answers. And by surveying approximately 10 % of the total truck fleet population each year, the survey is relatively cost challenging for Statistics Norway, binding up several FTEs both for data collection and data processing.

At the same time the field of road transport today is characterized by an increasing use of new technology, which in turn generate large amounts of data. There is widespread use of sensors on board vehicles for monitoring engine and car parts, driving behaviour and positioning. Also, road sensors generate continuous flows of data monitoring traffic throughout the road network. Lastly large amounts of data are also generated in enterprise resource planning systems, and logistics systems on the movement of types goods to-from destinations. In all these areas there might be a potential for utilizing new data sources for freight transport statistics, making data collection and statistical production more efficient and/or provide higher quality and precision. In this report we mainly investigate the area of

---

[22] The survey is also commonly named «survey on transport of goods». In this report we use the name used by Eurostat: https://ec.europa.eu/eurostat/documents/3859598/8918419/KS-GQ-17-114-EN-N.pdf/d9d20cec-d12c-491c-bb35-4fcf0ba6f9e0

sensor data from onboard trucks to get a better understanding of this landscape of data sources, and how this potentially can be utilized for producing official statistics.

## B.1 Business case for new data sources in freight transport statistics

The "big hunt for Big data" in the field of official statistics are present in various statistical areas, and has been for several years. Still, much of the initiatives are on an early stage, and many have also proven that it can be quite challenging before actual benefits to official statistics are materialised. When deciding upon whether one should investigate different fields of new data/big data one can take in two perspectives. The first is to what degree the existing statistics, and further data quality has a need for improvement, where both cost, quality and demand for more or new statistics should play a part. Secondly, potential alternative data sources must be prevalent, where technology, together with human and computer behaviour generates and stores data that might create a benefit for the statistical production. The "need for new data perspective" and "data potential perspective" should both be considered when deciding where to put resources in new/big data research.

As for the statistics on road transport, both of the mentioned perspectives are very much present. With obvious quality and cost challenges concerning today's data and data collection design, together with a huge landscape of sensor data generated in the freight transport industry, it makes the scope for looking for new data sources particularly interesting. This is not only the case in Norway, other EU-countries share much of the same prerequisites. Since the freight transport survey is mandatory across EU and EEA the information requirements are similar, and since the technology in trucks is more or less the same across Europe, This Norwegian study should of interest to other NSIs and Eurostat.

Some countries have already done their own investigations and testing to incorporate new data sources into road transport statistics, as a supplement to survey data. Statistics Netherlands (CBS) have come quite far by incorporating "weight in motion" (WIM) sensor data in the statistical production. This is devices installed beneath the highway to detect the axle weight of a truck passing. Such sensors are not used on Norwegian highways yet other than on a research/-testing scale, and thus will not be further investigated in this report.

## B.2 Today's Freight transport survey – overview of data and data collection design

The freight transport survey covers all trucks/tractive vehicles above 3,5 tons and collects data on both the movement of trucks and type of goods transported. Each quarter approximately 1 900 tractive vehicles are sampled from the Norwegian national vehicle register to participate in the survey (yearly sample of 7 000) (Statistics Norway 2019).

The enterprise who owns the sampled truck completes a diary questionnaire over a one-week period, filling in information on each trip the truck has carried out that week. Also trips where no goods are transported should be recorded. For each trip one must register the place of departure and the destination, the kilometres travelled, and some additional questions about the distance. Furthermore, there are several questions on what type of goods that was transported on each particular trip, with use of extensive dropdown lists for categorizing goods.

The information requirements from the Eurostat-regulation, commits the NSIs to provide microdata on all mandatory variables (For more info on variable requirements, see: (Eurostat 2019, chapter 10). The

microdata requirement explains why the survey is designed as it is, but it is also important when considering potential new types of data, and how they can be used together with other sources.

## B.3 Challenges with today's data collection design

The freight transport survey in Norway is conducted as a web survey, using SSB's standardized survey platform for business surveys. This platform is generic for all webform-based communication between businesses and government in Norway. However, it is poorly suited for web diary surveys and there are several questionnaire challenges and quality issues. As an example, the respondents have to repeat the same records manually many times for each trip registered, and there is little available assistance from automatically copying and reusing previous answers. Also, the technical solution is very rigid. The questionnaire must be filled out on a computer, preventing respondents, who often are the truck drivers, to use their phone or tablet for reporting right after or between trips.

The combination of a comprehensive questionnaire and a non-user-friendly survey solution, has a negative impact both on data quality and costs. Data quality is negatively affected by non-consistent fill-out of the diary questionnaire and a drop-out rate, with considerable resources being used on telephone follow-up by data collection staff. Although the official response rate is high (over 90 %) since the survey is mandatory, a longstanding problem has been a considerable over reporting of respondents who register the sampled truck as "out-of-service" in the survey period. Thereby they rule themselves out of the survey (almost 30 % of gross sample). This is not included in the official response rate, but some of it is a form of hidden non-response, as the actual share of trucks out-of-service must be lower. This has also been a challenge in other European countries.