# TUDelft

Delft University of Technology

## Big data of the past
## Analysis of historical freight shipping corridor data in the period 1662–1855

Wiegmans, Bart; Witte, Patrick; Janic, Milan; de Jong, Tom

**Important note**
To cite this publication, please use the final published version (if applicable).
Please check the document version above.

# Big data of the past: analysis of historical freight shipping corridor data in the period 1662-1855

Bart Wiegmans, Department of Transport and Planning, Faculty of Civil Engineering, TU Delft, the Netherlands, b.wiegmans@tudelft.nl

Patrick Witte, Department of Human Geography and Planning, Faculty of Geosciences, Utrecht University, the Netherlands, p.a.witte@uu.nl

Milan Janic, Department of Transport and Planning, Faculty of Civil Engineering, TU Delft, the Netherlands, m.janic@tudelft.nl

Tom de Jong, Stellenbosch University, Faculty of Economic and Management Services, Department of Logistics, tdejong@sun.ac.za

**Abstract**

This paper examines the use of big data and data analytics in international transport networks from the perspective of historical big data, focusing on shipping logs from the British, Dutch, Spanish and French fleets in between 1662 and 1855. Based on a large-scale database containing mainly meteorological data collected in the CLIWOC project (2003), we computed travel distances and analysed historical global maritime networks. This paper focuses on route choice, and consequently the time, distance, speed and reliability of the ships, covering different time periods, seasonal patterns and trade flows. The results reveal a clear picture of the main routes per nationality that is also indicative of the linguistical, cultural and economic colonial heritage that remains in the 'host' countries up to this day. The average daily distances covered vary over the countries involved, over the seasons and over different time periods. Also the trip characteristics vary notably over the different countries. Zooming in on the main trade flows, the corridor from the Netherlands to Indonesia stands out, but also considerable differences in average speed and stopover times were found along this route. Related to the complexity of using big data in studying international transport networks, our conclusion is that the degree of permutations and interactions with the dataset is not necessarily less for analyzing historical shipping records. It seems that big data of the past still can inspire future explorations of our historical transport networks on the world's oceans.

**Keywords**

data analysis, historical freight data, time, reliability, corridors, CLIWOC

# 1. Introduction

When looking at the use of big data and data analytics in international transport networks, the scientific focus in analyzing global maritime flows has been mainly on patterns of freight flows and analyses of the main transported freight (e.g. Ducruet, 2017). Ducruet et al. (2018) for instance relate the development of maritime networks to urban development and analyze long-term interdependencies between maritime networks and systems of cities, covering the period of 1890 until 2010. In these types of analyses, the use of empirical (big) data and data analytics is often instrumental, as is shown in a recent contribution by Wu et al. (2019) who conducted a vulnerability analysis of global container shipping liner networks based on disruptions at the important 'bottleneck' locations of the Malacca Strait, the Suez Canal and the Panama Canal. A similar approach is taken by Anchurra-Gonzalez et al. (2019) who have evaluated port disruption impacts in the global liner shipping network using game-theoretical models and cost-based container assignment models. A final recent example of using big data in transport is the port competitiveness model for ports along the Maritime Silk Road that was developed by Peng et al. (2018).

What is commonly visible in these studies – and in the popular debate around big data in general – is that big data are mostly consisting of data from numerous information-sensing Internet-of-Things devices such as large-scale user data from mobile devices, areal data (remote sensing), cameras, microphones, radio-frequency identification (RFID) readers and wireless sensor networks. Big data are mostly defined by their size, complexity and technologies, paying attention to the four V's of Volume, Velocity, Variety and Veracity (Ward & Barker, 2013; De Mauro et al., 2016). In contrast to looking mainly at the amount of Terabytes or the complexity of the technologies used (e.g. NoSQL), in this paper, we would like to stress the less commonly used definition of big data which mainly stresses complexity as a main factor, including a high degree of permutations and interactions with the dataset (Ward & Barker, 2013). In our view, this would also include working with historical data and records (see Wheeler et al., 2006 for insight into the complexity of analyzing historical records), which would open up possibilities for expanding the research scope of big data applications. In section 4, we give a detailed description and analysis of our big dataset regarding historical shipping.

Taking this as a point of departure, we observe in the literature different attributes that influence contemporary sea-route or land-mode choice (see Cullinane and Toy, 2000 for an overview). However, the more ancient periods, in particular the trade routes of Europe's colonial and imperial endeavors from the 17th century onwards, have – to our knowledge – never been analyzed in terms of route choice for maritime transports. However, from old shipping logbooks covering the period 1662-1855 and the countries France, the Netherlands, Spain, and United Kingdom (based on the open source shipping logbooks from the CLIWOC project, 2003) it shows that back then already time, speed, and distance were important variables to note. Additional route choice influencing attributes, but beyond the direct interest of the analyses presented in this paper, were weather conditions such as wind speed, wind force, and sight and sea conditions. Important other influences were the type of trades being made (such as cocoa, silver, gold, pepper, sugar, fruits, lumber, tobacco, flour, meat, furs, etc.) and their corresponding origins and destinations. Also, the prevailing wind patterns on earth played an important role in distance, speed and origin and destination selection. Also, safety was an important element at that time as the chances for a ship not returning home were considerable. The interest of this paper is looking at the consequences of wind patterns and types of trade in terms of route choice, and consequently the time, distance, speed and reliability of the Dutch, UK, French and Spanish fleets in general and specifically for different time periods, seasonal patterns and routes.

In particular, especially the time dimension for these global maritime freight flows has so far only received limited attention. In the paper by Slack et al. (2018), an analysis has been made of the ships time in port. Slack et al. (2018) analyzed the amount of time container vessels spend in port. They analyzed how the average vessel turnaround times (ATTs) vary among ports and they compared differences in ATTs with factors such as numbers of containers handled and several measures of port efficiency. They show that ATTs are differentiated regionally and functionally, rather than globally. This might also apply to the historical ships time and reliability (regional and functional differences between ships) that takes center stage in this paper. Distance performances might also be influenced for example by ship type and ship size, the transported freight, different trading routes, etc. In this paper, therefore, especially distance, time, speed and reliability (skewness and standard deviation) of historical sea sailing ships on different trade corridors are analyzed. Based on latitude and longitude data connected to daily patterns an average distance per ship and per day on maritime freight transport corridors can be calculated. Based on these distance data, we analyze historical freight transport patterns of the four different countries, further detailed to monthly periods, 25 year periods specifics and different trade routes. The focus is on historical time series and the patterns that can be derived from the database.
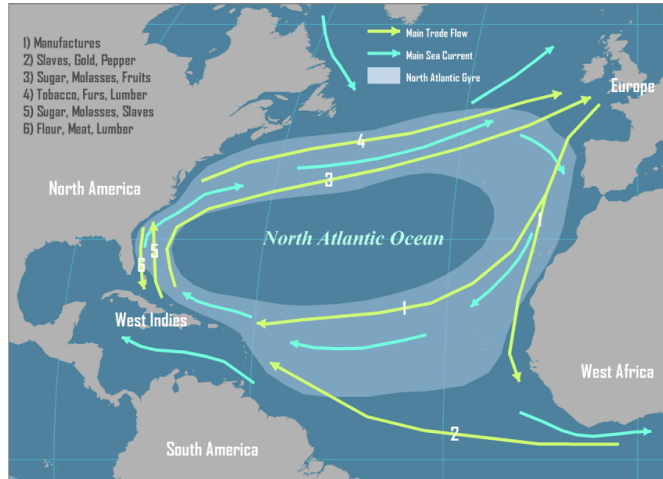
In the next section, the history of maritime freight transport corridors is discussed. Next, time and reliability of freight transport are defined from a historical perspective to be of practical value for use in this paper. Section four describes the dataset and the changes made to improve its quality. Section five contains the statistical analysis of the corridor levels for the Netherlands, UK, Spain, and France. Section six contains the conclusions of this paper.

## 2. The history of maritime freight transport and the main corridors

### 2.1 Historical development of maritime trade routes

During medieval times, the English Channel, the Baltic Sea, the Mediterranean Sea, and the North Sea were used for trade connecting important coastal and inland ports. Important goods traded at that time were wine, wool, salt, stone, grain, and timber and usually these were transported by galleys. In the 14th century, galleys were replaced by full-fledged sailing ships that were faster, could carry larger volumes and required smaller crews (such as the caravel, the carrack and then the galleon). At that time, already scale economies could be observed which stretches to today where still shipping companies seeks larger and faster ships that are cheaper to operate. In these centuries, the ships were not only trading ships but also carried cannons and other arms in order to protect the trades they were carrying. In the 1430s, the Portuguese discovered the North Atlantic circular wind pattern (also known as the trade winds) which signaled the start of the European expansion into the colonies (Figure 1). Comparable patterns were found on the Indian and Pacific oceans with the monsoon winds.

**Figure 1: Colonial Trade Pattern, North Atlantic, 18th Century**



Source: Rodrigue (2019)

In the 1450s, the traditional land trade routes connecting Asia with Europe were disrupted due to wars. This forced the European sea faring nations to find alternative maritime routes connecting the two continents. In 1492, following these efforts, Columbus discovered the American continent by sailing west. In 1497, Vasco de Gama sailed east and discovered a maritime route to India by rounding the Cape of Good Hope. Due to the possession of better armed, larger, and more efficient sailing ships the European powers (Spain, Portugal, France, United Kingdom, and the Netherlands) were able to control the seas (and thus international trade and colonization). The first private charter company emerged (The Dutch East India Company, or VOC), that established the first maritime trading network that spanned the entire world.

At this time, the hinterland transport system was still quite limited with most trade flows concentrating on port-port connections. By the late 16th century, inland canal systems started to emerge in Europe, initially in the Netherlands and in the United Kingdom. They enabled the large movements of bulk freight inland and facilitated the growth of regional trade. The Exeter Canal of 1566, designed to bypass a tricky stretch of river, was probably Britain's first dead-water canal with pound locks (Crompton, 2004). Crompton (2004) also pointed out that the fortunes of inland waterway transport have been more dependent on state support than those of other transport modes, even the railways. Maw et al. (2009) analyzed the Rochdale Canal, the busiest of Manchester's eastern waterways. This canal, with a length of 33 miles, was completed in 1804 for 600.000 Pounds and was used for the carriage of stone, coal, corn, textiles and cotton.
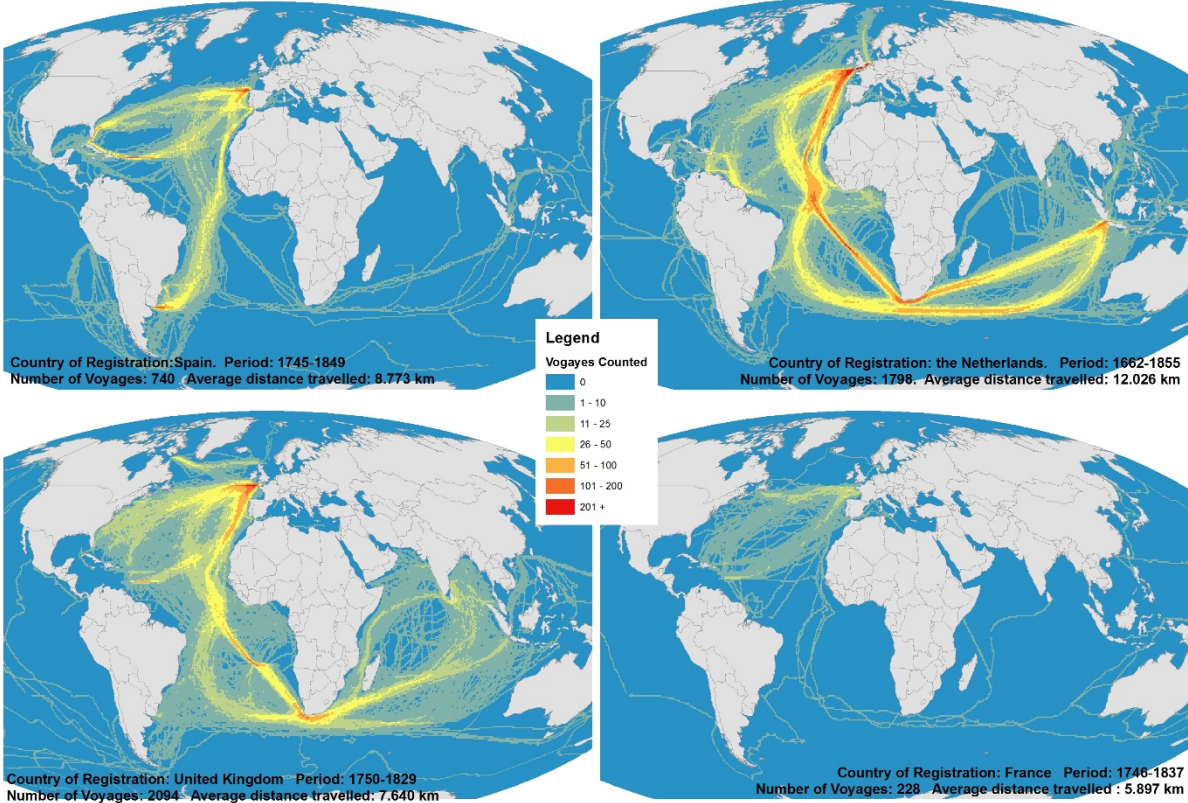
**2.2 Introduction to the main maritime corridors (1662-1855)**

It is interesting to analyze how the international trade corridors have developed since the opening up of both the 'west' and 'east' routes in the 15th century, and in particular since the European fleets started to use these routes for their colonial and imperial endeavors from the 17th century onwards. In its contemporary use, corridors can be viewed theoretically as bundles of infrastructure that connect two or more urban regions. These bundles usually exist in three modes: motorways, railway links and inland navigation and short sea connections (Priemus & Zonneveld, 2003; Chapman et al., 2003; Pain, 2011; Witte, 2012, 2014). However, in maritime history, sea corridors were the most important corridors and infrastructure was abundantly available. Sea transport tended to concentrate on certain origins and destinations (Figure 2).

From our own computations (Figure 2) we can see that the different countries involved (in our case Spain, the Netherlands, UK, and France) have followed their own international trade routes. We can see that Spain is mostly focused on the connections to and from the Caribbean, and to and from Latin America, oftentimes via the coast of Western Africa. Both the Dutch and UK fleets span the whole world, with both the 'western' and 'eastern' routes being covered. For the UK, there is a particular focus on the route to and from India, via the Cape of Good Hope. For the Dutch fleet there is a similar pattern, with the difference that beyond the Cape, the Dutch fleet mainly sailed for the Indonesian archipelago. For the French fleet, although data is more limited as compared to the other cases, especially the connection to and from Northern-America (east provinces of Canada) stand out. It goes without saying that all of these trade routes are not surprising when analyzing them in the light of the linguistical, cultural and economic colonial heritage remaining in the 'host' countries up to this day. After covering in the next section some working definitions for time and reliability that we aim to use in this paper, we continue to explain the data and empirical outcomes in more detail in the subsequent sections.

**Figure 2: Main maritime corridors following ship logbooks from 1662 to 1855 (own analysis)**



Ship log data from CLIWOC project: Main Sea Corridors by Country of Registration (Mollweide Projection)

Source: ship log data from CLIWOC project (2003). Data updated, computed and geocoded by authors (2019)

## 3. Time and Reliability in Maritime Freight Transport Corridors

For this paper, time and reliability in maritime freight transport corridors are two important factors to take into account. Therefore, in this section, we briefly present the definitions of time and reliability that we use for the empirical data analysis on the historical shipping records.

## 3.1 Time: transit times and delay

Time is a component quantity of various measurements used to sequence events, to compare the duration of events or the intervals between them, and to quantify rates of change of quantities in material reality or in the conscious experience. Current literature mostly focuses on transit times and delays (see e.g. Achurra-Gonzalez et al., 2019; Wu et al., 2019). Container carriers focus on liner services that have short transit times and high degrees of schedule reliability (Notteboom, 2006). Possible delays reduce the reliability of the transport service but might also incur additional logistics costs (such as additional inventory and production costs). Furthermore, vessel delays could also result in unproductive vessel times and even to the rescheduling of vessels. Especially the increase of the vessel speed can provide some 'slack' although at a cost (increased fuel usage).

Delays can be caused by four causes (Notteboom, 2006): terminal operations, port access, maritime passages, and chance. Increasing port volumes leading to capacity constraints can result in port congestion where vessels have to wait before being (un)loaded. Currently, this might be a more important issue than it was in the past when sail ships arrived in ports much less often and regularly. Delays in port access can result from lack of capacity for pilotage and towage services, sea lock capacity and tidal windows. Maritime passages such as the Panama Canal and the Suez Canal also influence carrier schedules in a significant way (see also Wu et al., 2019). Chance consists of weather circumstances, waiting times at bunkering sites, or mechanical problems. Shipping lines are constantly balancing factors such as risk of late arrivals, and the minimization of transit times. Port congestion is the main source of schedule unreliability (Notteboom, 2006).

In the case of the historical ships, time refers to days, as the historical shipping records often contain one measurement moment per 24 hours, where the exact geo-location (latitude, longitude) in combination with weather observations (wind force, wind strength, other weather conditions) were noted. Based on the time stamps in the shipping logs in combination with latitude and longitude indicators, we were able to calculate distances covered. Distance then refers to the average distances (converted to kilometers) covered per ship per day (24 hours). Averages were calculated by the total sailing distances per ship per trip, divided by the total number of days per trip as recorded in the shipping log files.

## 3.2 Reliability

Reliability, expressed as the variance of the lead time, is one of the key concerns of logistics decision makers (Dullaert and Zamparini, 2013). Reliability with regard to statistics refers to the overall consistency of a measure. Reliability in engineering refers to the ability of a system or component to perform its required functions under stated conditions for a specified time. For this paper, to avoid a semantic discussion over the definition of reliability (i.e. reliability of the ships themselves, or statistical reliability of the data), we propose a distinction between what we would call 'safety and security' and 'reliability'. With safety and security we mean the chances of ships actually completing the journeys (the shipwreck rate was considerable). It goes without saying that not many 'reliable' shipping logs of shipwrecks are preserved. With statistical reliability we mean the robustness of the general picture that emerges from the data (looking at basic indicators such as mean, standard deviation, range, skewness, Kurtosis), and whether journeys over time, over different corridors, over different seasons and over different countries show a more or less similar picture, and if not, how this can be explained.

## 4. Big Data and Methods

### 4.1 Introduction to the CLIWOC big dataset

A large-scale open source database exists that contains mostly meteorological observations made on board ships prior to 1856 (CLIWOC, 2003). So far, this source of software logs, as a historical forerunner of present-day big data, has mainly been used for meteorological purposes. The most significant logbook collections are from Spain, the UK, the Netherlands and France. Concerning the representativeness of this database (for a detailed account: see Wheeler et al., 2006) it can be stated that the Spanish data is largely complete (408 logbooks; 50.935 observations), the Dutch data is abstracted for about 50% of the total availability of logbooks (613 logbooks; 126.541 observations) and the UK data is more or less comparable to this in terms of volume (591 logbooks; 88.475 observations), but only covers a fraction of the total stock of British shipping logs. The biggest uncertainty is the French data, that accounts for only a rough 10% of the available logbooks (12 logbooks; 7.318 observations). Of course, statistically speaking, this cannot be justified as representative of the French situation. However, what we do intend to do, is to present the fullest picture as possible, given the data that we have. A few published logbooks from the USA, Germany, Denmark and Sweden were also included. Portugal, although a major player at the time, was not involved, because the logbooks could partly not be located, and partly could not be preserved. The resulting number of 273.269 observations that have been considered in the original CLIWOC project (2003) has in the meantime been updated to the most recent database of 287.114 observations (2019). This is the dataset we use for this paper.

The usage of this database enables the exploration of a long-existing source of reliable scientific information, namely ships' log books with its main focus on climate. Officers on board of these sail ships kept detailed logbooks of the ships' activities and management. Most of the logs are from vessels engaged in official government or military activities concerning trade, but a significant number have also survived from the ships of the quasi-governmental enterprises of British, Dutch and French trading companies. The data can be subjected to statistical analysis, used for synoptic reconstructions and for scientific interpretation and scrutiny. Initially, the database had a number of goals. First, to produce and make freely available for the scientific community the world's first daily oceanic climatological database for the period preceding 1850. Second, to realize the potential of the database to provide a better knowledge of oceanic climate variability over the study period. Third, to use the information to extend and enhance existing oceanic-climate databases. Fourth, to disseminate the project's findings and to stimulate interest and awareness in this source with a view to fostering its further development and realizing its scientific potential (Wheeler et al., 2006).

However, in our case, the usage of the database is extended to the analysis of historical maritime freight transport corridors that operated during these times. This database has so far not been exploited in detail with regard to large geographical ranges (i.e. trade corridors) and freight transport time and reliability. In this paper, we therefore specifically focus on analyzing differences between the countries involved (the Netherlands, United Kingdom, France and Spain) with regard to average daily distances, further detailed to monthly variance, developments over time periods of 25 years and differences in maritime trade corridors.

### 4.2 Characteristics of the big data in the database

Entries in the database are made chronologically (day by day) and include identification by geographic location (latitude, longitude), wind direction, wind force and other recorded weather elements. In selecting the data, the original CLIWOC researchers chose to focus on noon

observations, to only include voyages on open seas and oceans (thus excluding inland seas like the Baltics or Mediterranean) and to aim for as broad as possible temporal and spatial coverage, acknowledging the 'data gaps' due to political factors (e.g. the Napoleonic wars) and colonial factors (e.g. the lesser interest in the Pacific Ocean) (Garcia-Herrera et al., 2005). The database starts with the origin of the logbooks such as archive name, the country, and an ID. Next, the database contains ship characteristics such as ship origin, ship destination, ship name, ship type, and ship owning company. A third group of variables consists of Geodata (Table 1). This table can be used to find the longitude of the zero-meridian that was used in the logbooks. It also serves as lookup table for many of the used geographical names (anchor places, names of landmarks, ports sailed to and from, etc.). In the data calibration process we have improved the quality of the original dataset by adding and cross-checking geographic locations and by improving the match between latitude and longitude values and different spelling or names used in different languages (Dutch, English, French, Spanish) for ports and other destinations and points of interests.

**Table 1. Overview of Geodata**

| Field name | Description |
| --- | --- |
| Place | Original name (and spelling) of the place found in the logbook. Regularly more than one spelling or language were used for the same place. We decided to keep the original spelling intact, as much as possible, but to correct obvious errors (mostly misspellings or abbreviations) |
| Dec Latitude | Latitude in decimal degrees; North is positive, South is negative |
| Dec Longitude | Longitude in decimal degrees; East is positive, West is negative |
| Source | Name of the data source where the position of the place was found |
| Modern Name (i.e. English name) | Modern English name of the place |
| Alternative Name | Sometimes another name is common as well |
| Ocean | The name of the ocean where the place may be found |
| Spanish | Place name in Spanish |
| Dutch | Place name in Dutch |
| French | Place name in French |
| Notes | Additional notes |

Source: CLIWOC (2003)

A fourth group of variables consists of wind characteristics (wind force and wind direction) and the last group of variables consists of weather characteristics. Since the space in this paper is too limited to elaborate on variations in wind direction and wind force, and since this has already been done extensively by the original CLIWOC researchers, we have chosen to omit these variables from our empirical analyses. For a more detailed analysis of the climate results we refer to Koek and Können (2005).

**4.3 Big data management**

From this big data source we have made a selection of variables that enables us to focus on the geographical and transport aspects of the voyages. Included in the selection are: Voyage number, Record ID for the day, Voyage From (i.e. Origin), Voyage To (i.e. Destination), Ship Name, Nationality, Year, Month, Day, Latitude, Longitude, and added to the database are the Speed (in kilometres per hour) and the Distance (in kilometres per day). For (re-)calculating the distances, use has been made of the freeware GIS application FLOWMAP that can be used for analysis of flow data. In addition to this, an own script has been written to calculate the distances between chronological individual

records (using the Haversine formula). Several changes have been made to the original database to improve its quality and coverage. We encounter three main data management challenges, i.e. incomplete data, inconsistent data and how to handle modifications to the original data. Below, we describe how we dealt with this.

First, the dataset is not complete, meaning that incomplete data had to be deleted from the database. The following changes to the database have been implemented. In a first step, all records lacking a latitude or longitude have been deleted from the database. After that, all records were checked that were missing either the Voyage-ID, or the Origin, or the Destination. The next step involved deleting all records with no movement during subsequent days (same latitude/longitude). This can be observed often at the end of trips or when ships remain in ports for a number of days (for example between two different voyages). Limiting the database to the remaining cases, a column was added to this database with the number of days between two subsequent logbook entries. In most cases this is 1, sometimes 0 when more than one report is made per day (in case of multiple entries per day, the extra entries were omitted) and sometimes it is more than 1 (often this is between voyages, or when accidentally a record was found missing). In the final step, an additional column has been added that depicts the average latitude and longitude over 5 records and the deviation of one record from that average. All records with a deviation of larger than 5 degrees have been deleted from the database.

Second, the distances in the original database are not uniform (the dataset includes German sea miles, English sea miles, nautical miles, unknown units, etc.). Given the availability of the latitude and longitude coordinates, we have recalculated all the distances and put them in the same format. To this end, we used the before-mentioned FLOWMAP software and our own script. Based on latitude and longitude coordinates, we were able to calculate the distances and speed per day. Assuming a maximum speed of 20 knots per day that was typical for the most advanced ships in these days (approx. 37 km/hour) and assuming maximum efficiency in sailing (24 hours of continuous operation; 3 working shifts of 8 hours each), the maximum distance ships could theoretically cover is 888 kilometres per day. We used this distance to check and correct a few dozen of outlier cases. In many cases, these turned out to be wrongly entered or inconsistent records. These cases have also been removed from the database.

Third, we applied different data filters to run through the database in an explorative way, that allows us to understand the data patterns, including typical origins or destinations, or typical movements of ships during a particular season. This explorative analysis also revealed some additional odd cases. Examples include trips that were not recorded as a functional voyage from A to B, but rather a random exploration trip that was mostly recorded as 'coastal cruising' or using similar terminology. Since we want to focus on time and reliability of ships on their functional trading routes, we also chose to exclude the cruising/exploring trips from the data as much as possible.

All together, the original database of 287.114 records, has now been reduced to an improved database of 246.868 records with full data coverage for origin, destination, ship name, nationality, date of logbook entry, distance per day, speed per day, cumulative distance per trip, total days per trip, average trip length in days, average trip distance and average trip speed.

## 5. Results and analysis

The results section is structured in two parts. First, a descriptive account of the data is presented, starting with a general description, and then focusing on the daily distances, seasonal (i.e. monthly) patterns, long-term developments (periods of 25 years) and main trip characteristics for each of the

countries involved (UK, the Netherlands, Spain and France). Second, a more detailed account of the main trade routes is presented, zooming in on the most important trade corridor over the Atlantic Ocean, around the Cape and towards Indonesia, focusing mainly on the Dutch ships between the Netherlands and Indonesia.

## 5.1 Results of the descriptive analysis

### 5.1.1 General overview of the data

Table 2 shows the overall descriptive account of the records in the database (n=246.868). As can be observed, there is some other data next to the main data on the British, Dutch, Spanish and French ships. These few records have been entered by the Dutch project members of the CLIWOC project (Wheeler et al., 2006), but are not considered for further analysis in this paper because of the limited volume (less than 1.000 records in total and only 9 out of almost 5.000 trips). For the remaining data, it can be seen that especially the British and Dutch data have been collected in considerable volumes, given the number of records, number of trips and total distance covered. The Dutch ships show the highest average distances per day, the longest distances per trip and the longest durations per trip. These data are examined in further detail later (section 5.2). In the remainder of this section, the descriptive results for the British, Dutch, Spanish and French data is described.
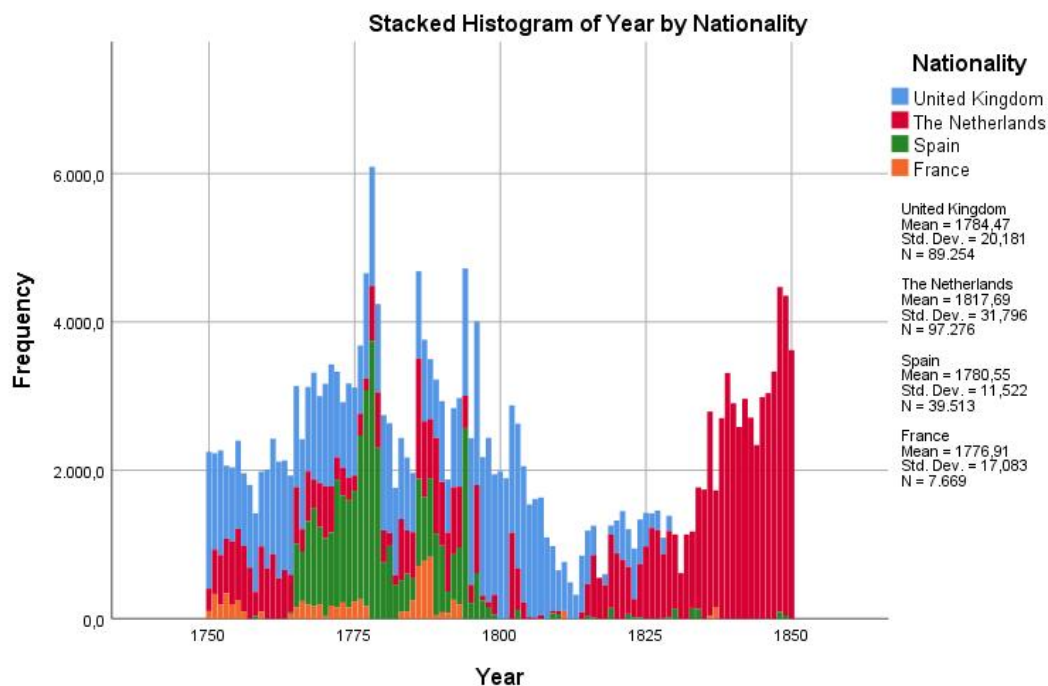
**Table 2: Descriptive overview of the dataset per nationality**

| Nationality | No. of records | No. of trips | Total distance in km | Avg. distance /day in km | Avg. distance /trip in km | Avg. days /trip |
|---|---|---|---|---|---|---|
| *America* | *202* | *3* | *35.334* | *173,34* | *11.778* | *67,3* |
| *Denmark* | *58* | *1* | *13.382* | *128,69* | *13.382* | *59,0* |
| *Germany* | *65* | *1* | *14.354* | *55,99* | *14.354* | *66,0* |
| *Sweden* | *631* | *4* | *119.276* | *113,83* | *29.819* | *237,5* |
| United Kingdom | 89.254 | 2.094 | 15.946.828 | 173,61 | 7.480 | 46,8 |
| The Netherlands | 109.276 | 1.799 | 21.622.095 | 195,14 | 11.749 | 63,3 |
| Spain | 39.579 | 737 | 6.433.740 | 158,08 | 8.469 | 56,3 |
| France | 7.803 | 229 | 1.334.989 | 167,33 | 5.816 | 36,2 |
| **Total** | **246.868** | **4.868** | **45.519.998** | **180,49** | **9.153** | **54,0** |

Source: authors' own adaptation from CLIWOC database (2003)

Zooming in on the data for the UK, Netherlands, Spain and France, we can see that most records are clustered in the period from the mid-eighteenth century until the mid-nineteenth century (1750-1850), although the range per country varies (Figure 3). The UK data range from 1750 to 1829, are more or less evenly distributed, but show a sharp drop in records after 1800. The Dutch data range from 1662 to 1855, with the main bulk of the data clustered between 1750 and 1850 and showing a sharp increase in records in the period of 1820-1850. This may have to do with the post-Napoleonic period and the increase in data collection afterwards. The Spanish data range from 1745 to 1849, with the bulk of the data clustered between 1765 and 1800. There is a stark decline in observations visible after 1780. Finally, the limited account of French data ranges from 1746 to 1837, with the bulk of the data clustered between 1750 and 1790. Most observations are to be found in the period from 1780 to 1800.

**Figure 3: Histogram of records per year by nationality**



N = number of sailing days
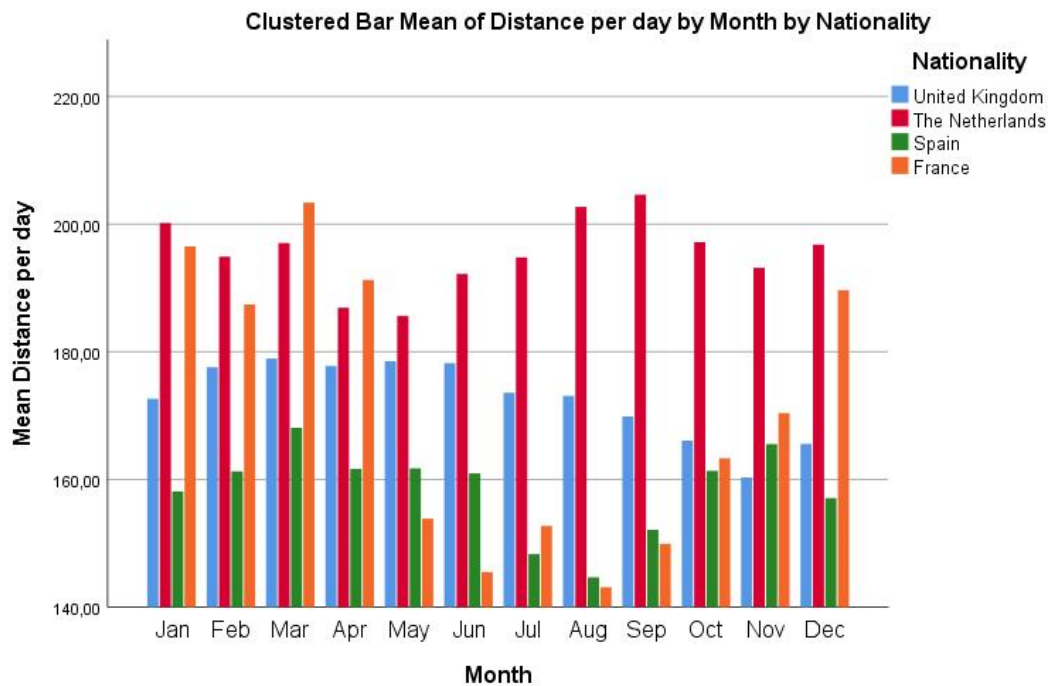Source: authors' own adaptation from CLIWOC database (2003)

*5.1.2 Daily distances*

When looking in more detail at the average total distance covered per day (180,49 km for all records), at a general level, we see some differences in the average distance per day, with the Dutch achieving the highest average distances (195,14 km) and the French the lowest (167,33 km). Although these differences are considerable, the standard deviations per country are not extremely different from the average standard deviation of 102,7 km, with the Dutch standard deviation being the highest (106,6 km) and the Spanish standard deviation being the lowest (91,3 km). The maximum distances per day are comparable (range 813,1 km for Spain to 839,0 km for the United Kingdom). This is all within the range of our assumed theoretical maximum distance of 888 kilometer per day (37km/hour). Finally, the overall distribution of records is slightly positively skewed, ranging from a Skewness value of ,259 for the Netherlands to ,645 for France (average: ,402) and a Kurtosis value of -,406 for the Netherlands versus ,676 for France (average: -,101). These results indicate a slight overrepresentation of lower than average daily distances, but not problematically so.

We can observe different patterns when comparing the data on a monthly or seasonal basis, and when looking at long-term differences over periods of 2,5 decades (25 years). First, when looking at the seasonal patterns (Figure 4), we can observe that the Dutch average daily distances per month fluctuate somewhat, with slight drops in April/May and October/November. Without being certain, this could indicate changes in speed owing to the changing monsoon winds in the in-between period (June to September), that usually result in relatively tranquil periods just before and after the changing of the winds, and the fact that most of the Dutch ships were travelling through territories where these monsoon winds are active (i.e. the Indian Ocean). Other explanations might relate to the types of trade these ships were undertaking, which might explain the increase in activities in the winter season for the French ships (e.g. trading North-American furs that are collected after the

hunting season). There are no variables in the dataset indicating the volume nor the type of goods being transported, so there is no way of being certain of these assumptions.
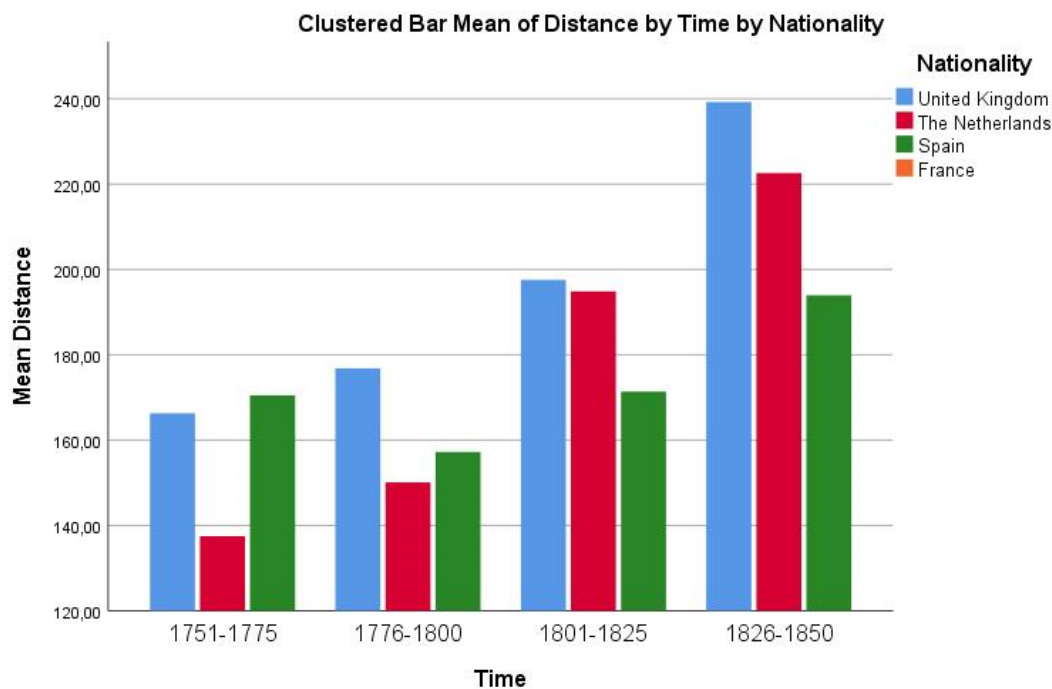
**Figure 4: Seasonal sailing patterns per nationality**



Source: authors' own adaptation of CLIWOC data (2003)

Second, taking a long-term perspective (Figure 5), we can tentatively see the efficiency of ships increasing. Especially the British and the Dutch ships show increasing average daily distances over the century-long period between 1750 and 1850. The UK average is increasing from 159,7 km in the period 1751 to 1775 to 237,7 km in the period 1826 to 1850. This last number should be taken with some caution, since it is only based on 895 records, compared to 31.170 records for the first period and also steam ships might be included. Still, the average of 194,2 for the 1801 to 1825 period based on 21.595 records is still a sign of increasing speeds and/or efficiency of operations. A similar pattern is observable for the Dutch data, showing an increase from 135,2 km (1751-1775; based on 15.081 records) to 220,4 km (1826-1850; based on 57.111 records). The picture for the Spanish data is more blurry, because of the before mentioned drop of records after 1800, leaving only about 1.000 records post-1800 compared to the roughly 38.000 records for the pre-1800 period. The French data is even more sketchy and incomplete, and is therefore left out of the figure below.

**Figure 5: Daily distance over different time periods by nationality**



Source: authors' own adaptation of CLIWOC data (2003)

*5.1.3 Main trip characteristics*

Another way of handling this vast database is by cumulating daily distances of individual shipping records per day into total trip distances and total trip durations (Table 3). This trims down the database from 246.868 individual (i.e. daily) records to 4.859 unique trips. A trip is usually recorded in the database as a single trip (e.g. Spain to Brazil). Doing so, new challenges emerge, and all countries involved have their own statistical challenges to deal with. In general, it can be stated that the variety increases. This is shown in the standard deviations for the British and Dutch data, that are almost as high or higher than the average trip distances and trip durations to which they refer. For instance, for the UK data, the mean trip distance is 7.480 km, with a standard deviation of 7.206 km. Looking at the trip duration, the standard deviation is even bigger than the mean and – the differences between minimum and maximum also being huge – the Kurtosis and Skewness values are exploding. The other countries show similar patterns, notably Spain and France. Nevertheless, the overall picture that emerges still is interesting for further analysis, and especially the Dutch data seem in general more or less reliable, given the data limitations that we have.
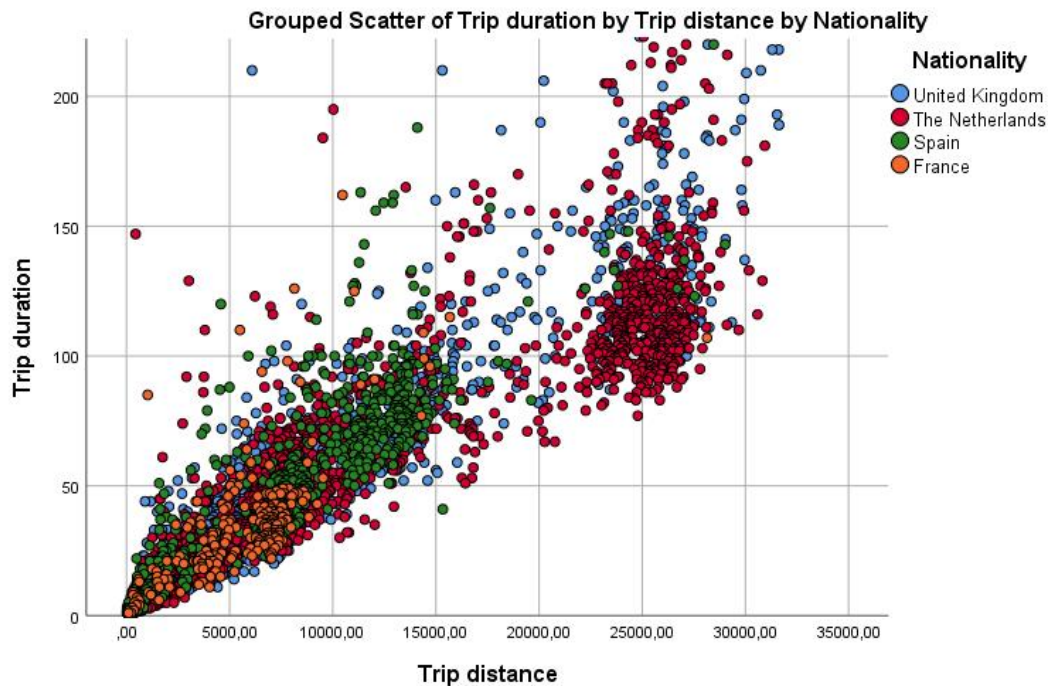
**Table 3: Trip distances and trip duration per nationality (n=4.859)**

|  | United Kingdom (n=2.094) | | The Netherlands (n=1.799) | | Spain (n=737) | | France (n=229) | |
|---|---|---|---|---|---|---|---|---|
|  | *Distance in km* | *Duration in days* | *Distance in km* | *Duration in days* | *Distance in km* | *Duration in days* | *Distance in km* | *Duration in days* |
| Mean | 7.480 | 46,8 | 11.749 | 63,3 | 8.470 | 56,3 | 5.816 | 36,2 |
| Std. Deviation | 7.206 | 53,4 | 9.749 | 50,2 | 4.978 | 41,1 | 3.100 | 23,4 |
| Minimum | 22 | 1 | 35 | 1 | 24 | 1 | 104 | 1 |
| Maximum | 40.976 | 1.165 | 34.505 | 369 | 33.358 | 458 | 28.135 | 162 |
| Kurtosis | 1,864 | 98,343 | -1,359 | 1,984 | 2,114 | 22,188 | 11,753 | 6,227 |
| Skewness | 1,575 | 6,215 | 0,472 | 1,045 | 0,678 | 3,200 | 1,778 | 2,072 |

n = number of trips. Source: authors' own adaptation from CLIWOC project (2003)

For further analysis, a scatter plot of trip distance versus trip duration is developed (Figure 6).

**Figure 6: Scatterplot of trip distance versus trip duration by nationality**



Source: authors' own adaptation from CLIWOC project (2003)

An interesting pattern emerges, with the French and British data mostly concerning smaller distances covered over shorter trips, Spain covering the intermediate range and the Dutch mostly covering large distances during longer trips. For interpreting this, we recalculated trip distance and trip duration in two separate, categorical variables. For trip distance, we use the categories short (<10.000 km), intermediate (10.000-20.000 km) and long (>20.000 km). For trip duration, we make a subdivision in four groups: <50 days, 50-100 days, 100-150 days and >150 days. This shows that, for the UK, three quarters of all trips (1.607 trips) are in the short distance range (<10.000 km) and within these trips nearly all trips (1.456) were conducted within 50 days. Of the countries involved, Spain has the highest share of trips in the middle distance range (34,5%), matching a comparable high share of trips in the two middle groups of the trip duration variable (224 trips with a 50 to 100 days duration). This may relate to the explorative voyages along the coast of Latin America and unto the Pacific Ocean. For the Netherlands, the distribution over the distance ranges is either short (58,4% of the trips) or long (30,5% of the trips). The middle range is hardly covered (11,2%). The high percentage of long distance trips may be explained by the explosive increase in long-distance trade voyages towards the Indonesian archipelago after Napoleon, in the first half of the nineteenth century (see Figure 3).

It is interesting to question why the difference between the short trips of the British ships and the long trips of the Dutch ships is so big, whereas they were roughly sailing the same routes and oceans (see Figure 2). This may relate to differences in keeping and recording the shipping logs by the naval officers, although there is no way of being certain just by examining the statistics. It could be that the UK data are recording different legs of the same voyage as different trips in our data (e.g. London-Accra / Accra-Cape Town / Cape Town-Calcutta), whereas the Dutch just record the origin and destination of the entire voyage (e.g. Rotterdam-Batavia), although in practice this was still cut into

different legs (e.g. Rotterdam-Accra / Accra-Cape Town / Cape Town-Batavia). To appreciate such differences, we need to dive into the data, by examining in more detail what is exactly happening along the main trade flows. It is to this matter that we turn now.

**5.2 Results of the GIS analyses of main trade flows**

*5.2.1 Top destinations and main trade flows*

Using geographical information systems (i.e. ArcMap) in combination with our own calculations, we were able to link the originally recorded origins and destinations of the voyages to the present day names of the ports as recorded in the World Port Index of the NGIA. In this way, 94% of the cases could be linked to a contemporary port location. Abstracting this to the country level, we can see the following patterns. First, a total of 417 unique origin-destination combinations on nations level could be found, of which 6 out of 10 of the most frequent combinations concerned three main routes of the British and Dutch fleets, including: Netherlands to Indonesia (253 trips), and vice versa (245 trips) often for spices trading, UK to Canada (97 trips), or vice versa often for wood and animal skins (74 trips) and Netherlands to Surinam (70 trips), or vice versa (65 trips) for cocoa, sugar and coffee. Second, looking at the top destinations of these two 'main users' of the world's oceans, we can see the following destinations standing out (Table 4).

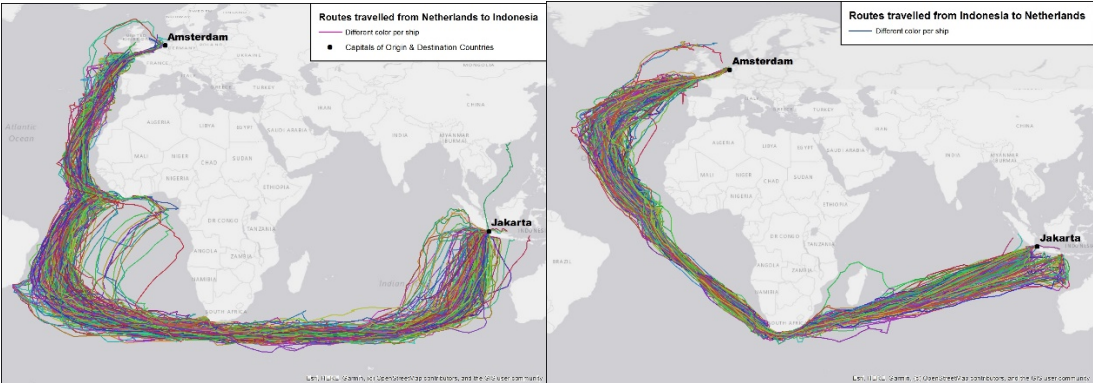**Table 4: Top destinations of trips for the British and Dutch data**

| The Netherlands | | United Kingdom | |
|---|---|---|---|
| *Destination* | *Number of trips* | *Destination* | *Number of trips* |
| Indonesia | 253 | Canada | 97 |
| Suriname | 70 | India | 69 |
| Ghana | 50 | Portugal | 59 |
| South Africa | 31 | Barbados | 44 |
| Svalbard | 24 | United States | 37 |
| Brazil | 16 | Spain | 29 |
| Netherlands Antilles | 14 | Indonesia | 23 |
| Antigua & Barbuda | 10 | South Africa | 19 |
| United States | 6 | Angola | 17 |

Source: authors' own adaptation from CLIWOC project (2003)

*5.2.2 A closer look at the Netherlands-Indonesia corridor*

Based on the top destinations, as a next step, we want to look in particular to the most frequently used route from the Netherlands to Indonesia, and vice versa. As can be seen in Figure 7, the routes for the home and return journey differ for the Netherlands-Indonesia route. For the home journey, following the so-called 'trade winds', ships would touch upon the West-African coast (in the Dutch case mainly stopping in between at Accra, Ghana) before sailing out almost to the Latin American coast (and sometimes disembarking in Rio de Janeiro, Brazil), before sailing to the Cape of Good Hope (South-Africa) and beyond. For the return journey, a more 'straight' path was chosen, sailing more directly from Indonesia to the Cape and, instead of touching either the Latin American or the West-African coast, going via the island of St. Helena in the southern parts of the Atlantic ocean (our data records 135 ships that have a stopover for two days or more at St. Helena), then towards the North-Atlantic ocean and 'bend' back again towards the port destinations of mostly Middelburg, Vlissingen, Rotterdam or Amsterdam in the Netherlands.
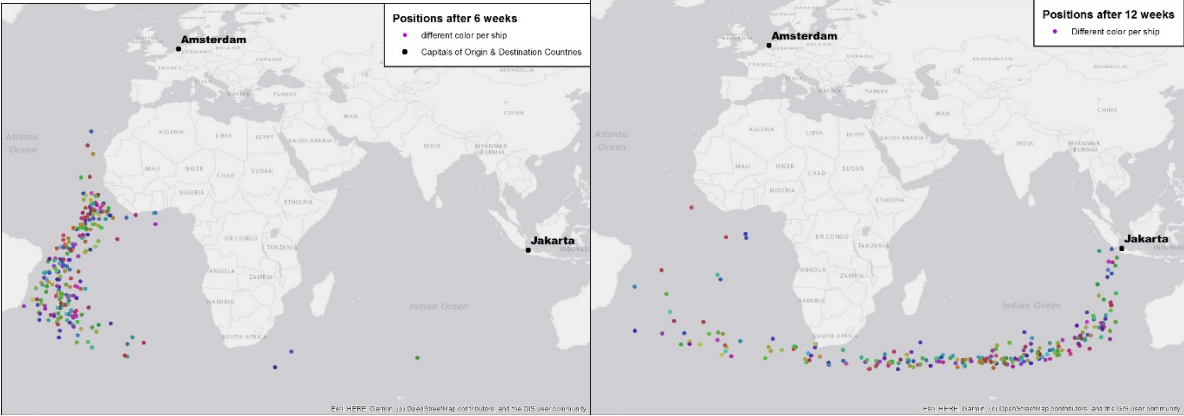
**Figure 7: Routes travelled from the Netherlands to Indonesia (left) and vice versa (right)**



*Routes: n=253                              n=245*
*Each coloured line represents a route covered by a ship*
Source: authors' own adaptation from CLIWOC project (2003)

Although on an aggregated and 'static' map this appears like a replicable recipe for success, we do see considerable differences when zooming in, especially regarding speed (Figure 8). <u>In the left part of the figure</u>, the progress after 6 weeks since the first recorded position is visualized, showing that the majority of the ships is still rather clustered – on their way from the West-African coast towards Latin America – while still some ships have barely passed by the Canary Islands (Moroccan coast) and some are already past the Cape. For this latter case, there may be some 'noise' in the data of ships cutting a voyage into two trips (e.g. Netherlands-Cape, Cape-Batavia), resulting in a considerable 'head start'. After 12 weeks (right part of the figure), the picture is even more dispersed, with a share of ships already at their destination, a great flock still on the Indian Ocean, and a considerable number still not yet beyond the Cape. After 14 weeks of travel since the first recorded position (112 days of sailing), 93 out of 253 ships have not yet arrived at their destination, and some ships have still more than 100 days of sailing ahead of them.

**Figure 8: Comparison of individual ships' progress after 6 and 12 weeks**
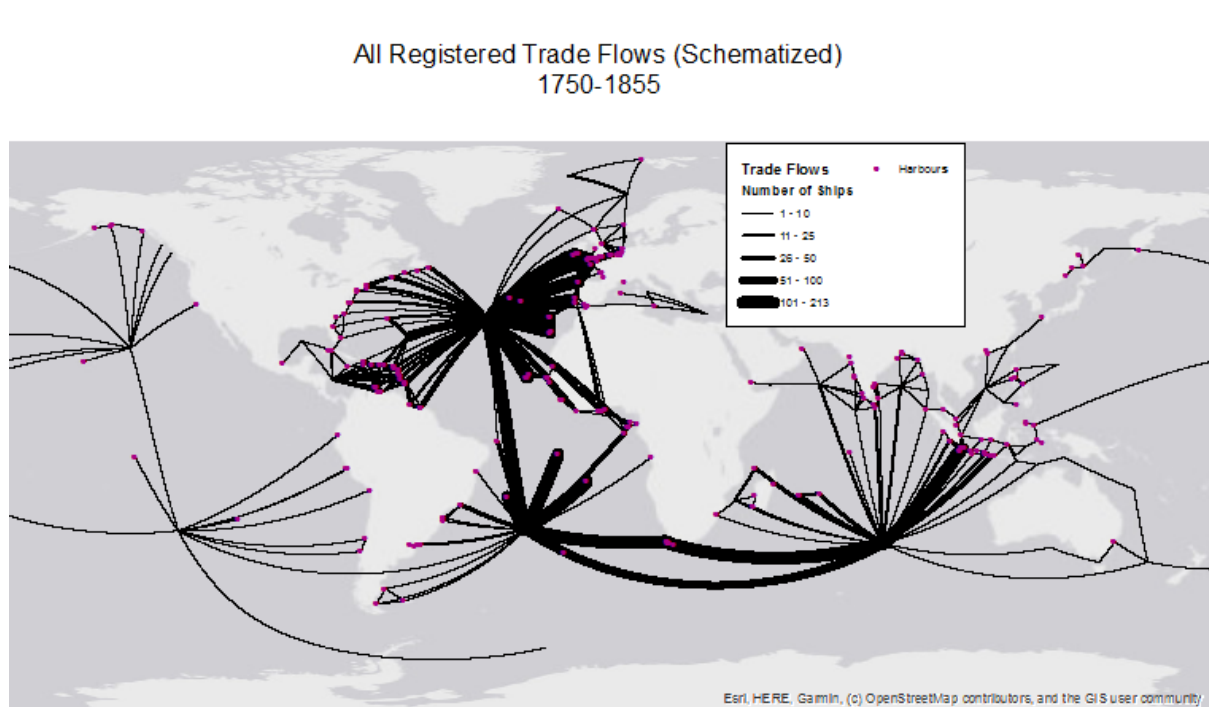


*Progress 6 weeks after first recorded position          Progress 12 weeks after first recorded position*
*The colors indicate ship positions after 6 and 12 weeks sailing respectively*
Source: authors' own adaptation from CLIWOC project (2003)

*5.2.3 From individual routes to schematized flows*

As a final step, we have developed a schematized overview of the worldwide trade flows (Figure 9). To do this, we applied the following procedure. First, we calculated the distance of the virtual connection between the first and last recorded coordinates of a trip (latitude, longitude) and the associated origin or destination according to the shipping log. We noticed sometimes large differences between the last recorded coordinates and the mentioned origin or destination. This is partly explained by incomplete data in the original entries; we assume that the first measurement of latitude and longitude was in some cases postponed until the ships were 'in the open' at sea. Another explanation might be that the CLIWOC project has just recorded the noon observation, whereas in practice there would perhaps have been multiple observations per day. A final explanation is that in the original 'data mining' and digitalization of the original shipping logs, the ports of origin and destination have been standardized and added as a separate extra coordinate at the beginning or end of a trip (e.g. in many cases the final coordinates of a trip would exactly match the coordinates of a certain port).

Second, to avoid these incompatible coordinates and port locations, we selected the virtual connections that were within a 350 km radius, leaving 1.250 trips for analysis, with an average virtual distance between the last coordinates and the ports of 144 km. The selected cases start from or terminate in 235 different ports. Third, we aggregated all coordinates within the 350 km range from the port, and in a similar way aggregated the other coordinates to certain segments of the worldwide oceans and waters. Fourth, all resulting dots were connected chronologically and counted for the number of ships on every connection. The result is depicted in Figure 9, clearly showing the main Atlantic backbone, and the subsequent voyages beyond the Cape, either to India or to Indonesia.

**Figure 9: Schematized trade flows based on the shipping logs (n=1.250 trips)**
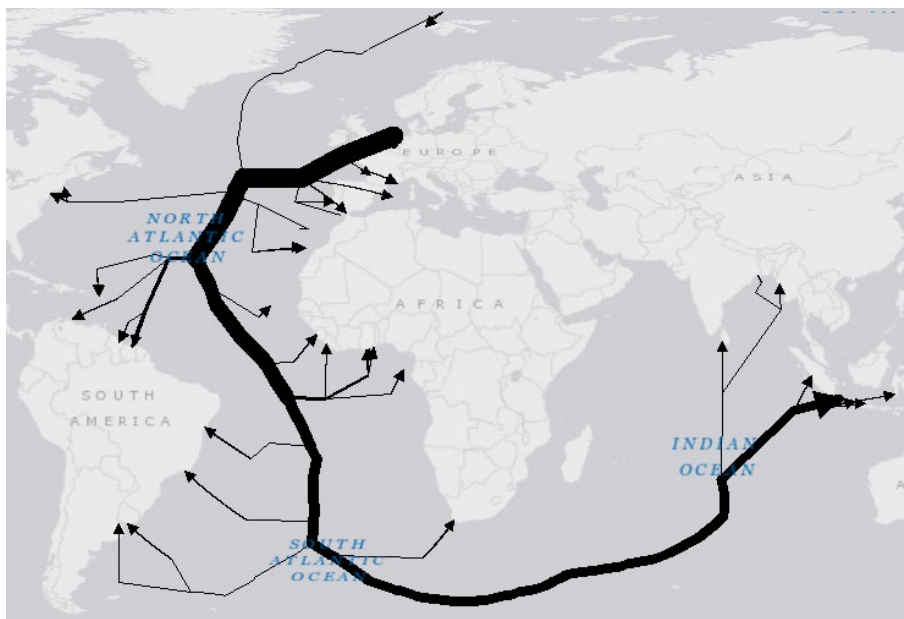


Source: authors' own adaptation from CLIWOC project (2003)

When the schematized trade flows network is further detailed to an artificial (spatial) neural network of connected origins and destinations of Dutch vessels, the following picture emerges (Figure 10). The

map shows the network of Dutch vessels also having a Dutch port of origin in the period 1750-1850. This results in 493 ships in total, which means roughly 5 ships per year. The thickness of the arrows is proportional to the number of ships on that part of the route. Of the 43 different destinations of these ships, the connection with Indonesia once again stands out, with 213 trips to Batavia (nowadays known as Jakarta). Other important destinations are Paramaribo in Surinam (64 trips), Tema in Ghana (port location close to Accra, 37 trips), Ny-Alesund on Svalbard (22 trips), Simonstown in South Africa (18 trips), Rio de Janeiro in Brazil (18 trips) and Willemstad on Curacao (Netherlands Antilles, 13 trips). So, Jakarta (Indonesia) is by far the most popular destination, with an average of two ships per year.

**Figure 10: Spatial neural network map of Dutch vessels from 1750 to 1850 (n=493)**



Source: authors' own adaptation from CLIWOC project (2003)

**6. Conclusions**

This paper has looked into the use of big data and data analytics in international transport networks from the specific perspective of historical big data, more specifically shipping logs from the British, Dutch, Spanish and French fleets in between 1662 and 1855. As an historical antecedent of present day big data, the roughly 300.000 shipping records from the CLIWOC project (2003) provide a complex database for analyzing historical global maritime freight flows. We have specifically focused on route choice, that was likely based on wind patterns and type of trade, and consequently the time, distance, speed and reliability of the ships in our database in general, and specifically for different time periods, seasonal patterns and trade flows.

The results reveal a clear picture of the main routes per nationality that is also indicative of the linguistical, cultural and economic colonial heritage that remains in the 'host' countries up to this day, with the French ships focusing mostly on Canada, the Spanish ships on Latin America, the UK fleet having a global focus with North-America and India as main focal points, and the Netherlands having clear 'east' and 'west' routes to the Indonesian archipelago and Surinam and the Netherlands Antilles respectively. The total recorded distance of nearly 46 million kilometers (almost 600 times around the globe) spread over nearly 5.000 unique trips equals an average daily distance of 180 kilometers. This is subsequently varying over the countries involved (with the Dutch ships achieving

the highest averages), over various seasons (with the highest daily distances in the monsoon winds season) and over longer time periods (with increasing daily distances over time, possibly indicating technological and navigational innovations of the ships).

Looking at trip characteristics (distance and duration), we have found different behaviors for the different nationalities involved. The Dutch fleet was mainly concerned with the long-distance voyages taking several months to complete, whereas the British ships were mainly active in the short distance ranges and taking a relatively short time to complete trips. The Spanish activities were mainly clustered in the middle range, mostly indicating explorative voyages along the Latin American coast and onto the Pacific Ocean. Zooming in on the main trade flows, the corridor from the Netherlands to Indonesia via the Atlantic backbone and the Cape of Good hope stands out. In the period from 1750 to 1850 there were regular connections between the Dutch ports and Batavia (present day Jakarta), but we also noticed considerable differences in average speed and stopover times along this route, with the island of St. Helena being the most frequent stopping place.

Related to the complexity of using big data in this type of research, our conclusion also is that the degree of permutations and interactions with the dataset is not necessarily less for analyzing historical shipping records. Although this dataset is 'just' 300.000 records and its size is nowhere near Terabytes of data, the amount of in-depth, 'manual' labor to reach satisfactory levels of reliability is considerable, as is shown in the analyses above. Similar to present day big data of for instance smartphone GPS locations, we also face problems of, for instance, how to deal with unlikely outliers and how to investigate what has exactly happened along a certain route. With more detailed and advanced types of analyses, the number of suitable and reliable records was dropping notably.

This does not mean that nothing can be learned from our findings. The managerial implications of this paper can for instance be related to current day practice in deep-sea shipping. The data analytics performed on this database shed light on historical time (distance) and reliability figures. With the increasing efforts in present day deep-sea shipping to become more sustainable, adding sails to existing ships is one of the potential successful interventions. Historical travel patterns can be used for future routing of combined diesel/sailing propulsion, to optimize the use of the wind and reduce diesel usage as much as possible. The relevance of this suggestion is also shown through the routing of the current Volvo Ocean Races, that still use largely the same routes and sailing patterns as are indicated in our flow maps.

Further research can focus on more detailed elements of the database, such as detailing certain routes and analyzing in-depth the stopping places and activities along these routes. For instance, for the Dutch data, the trade routes of the West and East Indies trading companies to Surinam and/or the Dutch Antilles and Indonesia could be further explored. The current dataset does not provide many insights in types of trades and details of activities and stops along the way. However, a first exploration of the archives of the East Indies Trading Company (VOC) website and the individual shipping journals that are archived in the EASY dataset of the DANS project (Data Archiving and Networked Services) shows potential for more elaborated types of analysis. This might shed more light on issues such as transported freight, freight values, ship owners, risks and mitigation strategies, military operations and the like. Further research could also consist of analyzing the dynamic changes of navigation characteristics in different countries and regions over such a long time period. In addition, if the characteristics of navigation after 1885 could be connected with the historical results in this article, even more interesting results may be obtained. Possible big data sources to make this connection include the US Maury database, or the ICOADS database. All in all, it seems that big data of the past still can inspire future explorations of our historical transport networks on the world's oceans.

# References

Achurra-Gonzalez, P.E., Angeloudis, P., Goldbeck, N., Graham, D.J., Zavitsas, K. & M.E.J. Stettler (2019), Evaluation of port disruption impacts in the global liner shipping network. Journal of Shipping and Trade 4 (3), pp. 1-21.

Chapman, D., Pratt, D., Larkham, P., Dickins, I., 2003. Concepts and definitions of corridors: Evidence from England's Midlands. Journal of Transport Geography 11(3), 179-191.

CLIWOC (2003) Ship log data from CLIWOC Project (online, accessed January 2019), http://projects.knmi.nl/cliwoc/download/CLIWOC21lim.htm.

Crompton G. (2004) The tortoise and the economy; inland waterway navigation in international economic history, The Journal of Transport History, 25(2), 1-22. DOI: 10.7227/TJTH.25.2.1.

Cullinane K. Toy N. (2000) Identifying influential attributes in freight route/mode choice decisions: a content analysis, Transportation Research Part E: Logistics and Transportation Review 36 (1), 41–53.

De Mauro, A., Greco, M. & M. Grimaldi (2016) A formal definition of Big Data based on its essential features. Library Review 65:3, p. 122.

Ducruet C. 2017 Multilayer dynamics of complex spatial networks: The case of global maritime flows (1977–2008), Journal of Transport Geography, 60, 47-58.

Ducruet, C., Cuyala, S. & A.E. Hosni (2018), Maritime networks as systems of cities: The long-term interdependencies between global shipping flows and urban development (1890–2010). Journal of Transport Geography, 66, 340-355.

Dullaert W. Zamparini L. 2013 The impact of lead time reliability in freight transport: a logistics assessment of transport economics findings, Transportation Research Part E, 49, 190-200. http://dx.doi.org/10.1016/j.tre.2012.08.005.

García-Herrera R. Wilkinson C. Koek B. Prieto M.R. Clavo N. Hernández, E. (2005) Description and general background to ships' logbook as a source of climatic data, Climatic Change, 73, 13-36. DOI: 10.1007/s10584-005-6954-4.

Koek, F. & G. Können (2005) Determination of wind force and present weather terms: the Dutch case. Climatic Change 73, 79-95.

Maw P. Wyke T. Kidd A (2009) Water transport in the Industrial Age: Commodities and carriers on the Rochdale Canal, 1804-1855, Journal of Transport History, 30 (2), 200-228. http://dx.doi.org/10.7227/TJTH.30.2.6.

Notteboom, T.E. (2006) Time factor in liner shipping services, Maritime Economics and Logistics, 8, 19-39. https://doi.org/10.1057/palgrave.mel.9100148.

Pain, K., 2011. 'New Worlds' for 'Old'? Twenty-first-century gateways and corridors: reflections on a European spatial perspective. Int. J. Urban Reg. Res. 35 (6), 1154–1174.

Peng, P., Yang, Y., Lu, F., Cheng, S., Mou, N. & Yang, R. (2018), Modelling the competitiveness of the ports along the Maritime Silk Road with big data. Transportation Research Part A, 118, 852-867.

Priemus, H., Zonneveld, W., 2003. What are the corridors and what are the issues? Introduction to special issue: the governance of corridors. Journal of Transport Geography 11, 167-177.

Rodrigue, J.P. (2019) Colonial Trade Pattern, North Atlantic, 18th Century (online, accessed January 2019), https://transportgeography.org/?page_id=1094.

Slack, B., Comtois, C. Wiegmans, B. Witte, P.A. (2018) Ships Time in Port, International Journal of Shipping and Transport Logistics, 10(1), 45-62.

Ward, J.S. & A. Barker (2013) Undefined By Data: A Survey of Big Data Definitions. arXiv:1309.5821v1.

Wheeler, D., Garcia-Herrera, R., Koek, F., Wilkinson, C., Konnen, G., del Rosaria Prieto, M., Jones, P. & R. Casale (2006) CLIWOC: Climatological Database for the World's Oceans: 1750 to 1850. Results of a research project EVK1-CT-2000-00090. Brussels: European Commission.

Witte, P., van Oort, F., Wiegmans, B., Spit, T. (2014). European corridors as carriers of dynamic agglomeration externalities?, European Planning Studies, 22(11), 2326-2350. DOI:10.1080/09654313.2013.837153.

Witte, P., Wiegmans, B., Oort, F.G. van, Spit, T.J.M., (2012). Chokepoints in corridors: Perspectives on bottlenecks in the European transport network. Research in Transportation Business and Management, 5, pp. 57-66.

Wu, D., Wang, N., Yu, A. & Wu, N. (2019), Vulnerability analysis of global container shipping liner network based on main channel disruption. Maritime Policy & Management, 46 (4), 394-409.