

## RESEARCH ARTICLE

# Towards the Internet of Water: Using Graph Databases for Hydrological Analysis on the Flemish River System

Erik Bollen<sup>1</sup>, Rik Hendrix<sup>2</sup>, Bart Kuijpers<sup>3</sup>, Alejandro Vaisman<sup>4</sup>

*(Received 00 Month 200x; final version received 00 Month 200x)*

**Keywords:** Graph Databases, Graph OLAP, Graph Aggregation, Hydrology

The “Internet of Water” project will deploy 2500 sensors along the Flemish river system, in Belgium. These sensors will be part of a monitoring system. This will produce an enormous amount of data, over which prediction and analysis tasks can be performed. For representing, storing, and querying river data, relational databases are normally used. However, this choice introduces an “impedance mismatch” between the conceptual representation (typically a graph) and the storage model (relational tables). To solve this problem, this paper proposes to use graph databases. The Flemish river system is presented as a use case and the Neo4j graph database and its high-level query language, Cypher, are used for storing and querying the data, respectively. A relational alternative is implemented over the PostgreSQL database. A collection of representative queries of interest for hydrologists is defined over both database implementations.

---

<sup>1</sup>Hasselt University, Data Science Institute and Database and Theoretical Computer Sciences Research Group, Belgium and VITO, Flemish Institute for Technological Research, Mol, Belgium; email: erik.bollen@uhasselt.be

<sup>2</sup>VITO, Flemish Institute for Technological Research, Mol, Belgium; email: rik.hendrix@vito.be

<sup>3</sup>Hasselt University, Data Science Institute and Database and Theoretical Computer Sciences Research Group, Belgium; email: bart.kuijpers@uhasselt.be

<sup>4</sup>Instituto Tecnológico de Buenos Aires, Buenos Aires, Argentina; email: avaisman@itba.edu.ar

## 1. Introduction and Motivation

Recent climate changes are increasingly leading to extreme meteorological weather phenomena. These situations have impact on water supply and water quality, for example, due to the influence of the salty sea on rivers, which can have a negative impact on the water and surrounding land area. Monitoring water quality and quantity is becoming more and more relevant. In Belgium, the “Internet of Water” project<sup>1</sup> (IoW) aims at enhancing monitoring and governance of the Flemish waterways. This project plans to deploy 2,500 sensors along the Flemish river system. These will allow, for example, to trigger a warning if certain measurements pass over pre-defined thresholds. A second example is, if a pollution problem is detected by a sensor at a certain location, the state at downstream locations could be predicted in order to take timely appropriate action. Also, typical data analysis tasks can be performed using the enormous amount of data that will be produced. All of the above requires appropriate modelling, storing and querying of such data. Normally, this would be done using relational databases. Nowadays, graph databases also appear as good candidates for these tasks, as the following discussion suggests.

Property graphs (Robinson *et al.* 2013), that is, graphs whose nodes and edges are annotated with properties, are typically used to model networks (e.g., social networks, sensor networks) to perform data analysis. The property graph data model is an abstraction that can also be used to represent rivers in a natural way. For example, using this model, the river segments can be represented as nodes in a graph, and an edge would go from one segment to another, if they are consecutive in the direction of the flow. In addition, spatio-temporal coordinates can be included as properties, as well as other characteristics of the river segments. Also, hierarchical contextual data could be defined, which would allow representing the graph at different granularities, for analytical querying involving data summarisation. Modelling rivers using graphs allows storing them in a natural way, using graph databases (Angles 2012, 2018), rather than relational databases, preventing the “impedance mismatch” problem, that arises when the natural network structure is split into many records of a relational table. For example, when a river network is stored as a graph, and represented as indicated above, finding a path between river segments is straightforward using a native graph database<sup>2</sup>. On the other hand, using a relational database, segments would be represented as rows in a table, therefore, finding a path requires self-joining the table as many times as the length of the path requires. In particular, in this paper, the Neo4j graph database<sup>3</sup> is used. Besides its popularity, the Neo4j community has developed several libraries of functions, that are easily added to the database as plugins. These libraries include a powerful machinery of algorithms for finding paths in graphs, handling many different data types and performing usual data science tasks. There is also a spatial library<sup>4</sup>, which can enhance the analysis possibilities. Last, but not least, Neo4j comes with a high-level graph query language, Cypher.

This paper proposes the use of graph databases for facilitating the work of hydrologists along two main dimensions: On the one hand, certain queries of interest can be expressed intuitively by non-expert professionals; On the other hand, more involved queries may,

---

<sup>1</sup><https://www.internetofwater.be/en/what-is-internet-or-water/>

<sup>2</sup>Graph databases are called native if they use specialised data structures for storing data. On the contrary, if they provide interface for other kind of storage, e.g., relational databases, they are called non-native.

<sup>3</sup><http://www.neo4j.com>

<sup>4</sup><https://github.com/neo4j-contrib/spatial-algorithms/releases/tag/0.2.3-neo4j-4.1.3>

sometimes, run faster over the graph database than over the relational alternative, thanks to specialised native data structures that allow efficient path traversal. Concretely, the work tackles the following questions: (a) *Can graph databases be successfully used to model, store, and query river flows?* (b) *If so, which are the kinds of queries that could benefit the most from this approach?* (c) *Is it more intuitive and simple for a non-expert user, to express queries using high-level graph query languages than writing SQL queries over a relational database?* To answer this questions, the Flemish river system is studied and discussed in depth. Further, the process of transforming the source data into a format suitable for querying is also addressed in this work.

In summary, the contributions of this paper are:

- (1) The definition of a property graph data model for representing river systems, that can be extended to other kinds of transportation networks.
- (2) A real-world case study of this proposal, using the complete Flanders river system.
- (3) A description of the data acquisition and transformation processes, that take the river system data from a shapefile into a relational database, create a graph, and store this graph using graph databases.
- (4) A definition and analysis of a collection of queries, expressed in Cypher and SQL, and executed over the Neo4j and the PostgreSQL databases, respectively. The queries are run and the results discussed and reported.

It follows, from the experiments and the analysis that, in most of the cases, queries over the graph database show better performance (with a few exceptions) than their relational equivalent, particularly in the queries asking for paths. Also, in many cases, queries are more easily and naturally expressed in Cypher than in SQL. However, for some queries, good performance is achieved at the expense of writing more complex Cypher expressions, which are not very intuitive.

The remainder of this paper is organised as follows: in Section 2 related work is discussed. The problem of acquiring and preparing the river data is discussed in Section 3, and in Section 4, the relational and graph storage are described and discussed. A case study is presented in Section 5 where a collection of queries, to analyse the data in relational and graph databases, are proposed. An experimental evaluation of these queries, over Neo4j and PostgreSQL, is reported and discussed in Section 6. Finally, Section 7 concludes the paper, also addressing future work and open problems.

## 2. Related Work

This section studies related work, starting from a description of the context of the problem, namely the rivers in the Flanders region in Belgium. Then, graph databases are discussed. Finally, a brief comparison between relational and graph databases is presented.

### 2.1. Data-driven approaches for studying flows in river systems

The region of Flanders is located in the northern part of Belgium. In spite of encompassing a relatively small area, watersheds within Flanders exhibit a wide range of regimes which require localised parameterisations, for more accurate hydrological modelling (Heuvelmans *et al.* 2004). In recent decades, the chance of extreme meteorological events has increased in Belgium. This includes the occurrence of heavy storms and fre-

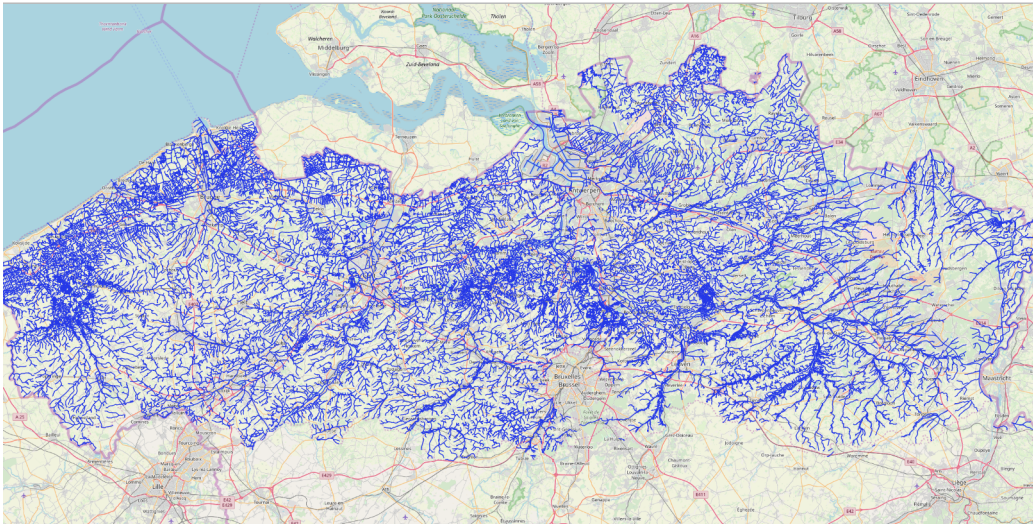


Figure 1. Overview of the Flemish river system.

quent heatwaves, which results in increased floods and drought periods (Brouwers *et al.* 2015). Drier periods, specifically, have a dual impact in the region, as less rainwater runoff causes higher risks of seawater intrusion from the North Sea resulting in salinisation of groundwater and soils. As more than 50% of the area in Flanders is used for agriculture, such events severely impact the country's socio-economic status (Gobin 2012). The implementation of a dense sensor network over a hydrologically complicated and environmentally vulnerable region, allows building an integrated geospatial data-driven river system. To put the problem in context, Figure 1 depicts an overview of the Flemish river system, using QGIS<sup>1</sup> over an OpenStreetMap background. The vast amount of river branches can be clearly seen in Figure 1.

Physical process-based modelling described above, although necessary, does not suffice for the current vast amounts of data from various sources for real-time applications. Additionally, commonly used spatially distributed hydrologic models still rely, to some extent, on empirical parameterisations and extensive calibration. The implementation of complementary data-driven approaches have become increasingly popular and have successfully represented hydrological processes (Ahani *et al.* 2018, Solomatine and Ostfeld 2008). Data-driven approaches allow for additional insights based on classifications or clustering of regions with similar input and output relationships at varying spatial or temporal resolutions, a task not easily implementable in traditional process-based models.

Typically, relational databases are considered as the standard systems for storing data like the one needed for the study introduced above. However, as argued in Section 1, graph databases appear as natural candidates for that task, since the river topology can be modelled as a graph, and stored in native data structures, appropriate to answer the required queries efficiently. (Demir and Szczepanek 2017) extensively discuss graph data models as a natural way of representing river networks. In fact, they simplify the analysis, arguing that a *tree* representation would suffice to cover a hydrologist interest. The present paper shows (see Section 5, Query 5.7) that this simplification is not very realistic. Also, benchmarking was performed over a small tree containing one thousand nodes, and using PostgreSQL to store the graph data model. That is, the approach does

<sup>1</sup><https://www.qgis.org>

not tackle the discrepancy between the conceptual and physical data models. The work by (Daltio and Medeiros 2015) addresses this issue, proposing using Neo4j for storing a river network. They present “Hydrograph”, a tool to explore geographic data in graph databases. The graph data model adopted in this work is a binary tree. However, the work does not report implementation details, or queries over the river system.

We note that both cases discussed above simplify the problem adopting a tree representation. Opposite to this, in the present paper an actual graph representation is assumed, which, as will be discussed below, is much more demanding, since it increases the number of possible paths, but, on the other hand, realistically represents the river flow. In addition, no tests on graph databases are reported in such efforts. In this paper, relational databases (in this case, PostgreSQL) and graph databases (Neo4j) are tested under a collection of typical queries required in river analysis and parameter prediction.

Relational database technology is mature and well-known, while graph databases are relatively new, therefore a brief description is provided next.

## 2.2. Graph Databases

In the context of graph databases, two models are used in practice:

- (a) Models based on RDF<sup>1</sup>, oriented to the Semantic Web; and
- (b) Models based on property graphs.

Models of type (a) represent data as sets of triples where each triple consists of three elements that are referred to as the subject, the predicate, and the object of the triple. These triples allow describing arbitrary objects in terms of their attributes and their relationships to other objects. Informally, a collection of RDF triples is an RDF graph. In models of type (b) (Angles *et al.* 2017), nodes and edges are labelled with a sequence of attribute-value pairs. It is an extension of classical graph database models, frequently used for implementations in practical applications. The main reason for storing attributes in nodes and edges is speeding up the retrieval of the data directly related to a certain node. For an extensive and comprehensive bibliography on graph database models, the interested reader is referred to (Angles and Gutierrez 2008, Angles 2018). Although the models of type (a) have a general scope, the structure of RDF makes them not as efficient as the other models, which are aimed at reaching a local scope. An important feature of RDF-base graph models, however, is that they follow a standard, which is not yet the case for the other graph databases, therefore they are typically used for metadata representation. Many works have proposed RDF to annotate trajectories with semantic information (Fileto *et al.* 2015, da Silva *et al.* 2015, Ruback *et al.* 2016). Hartig (Hartig 2014) proposes a formal way of reconciling both models formally, through a collection of well-defined transformations between property graphs and RDF graphs. He shows that property graphs could, in the end, be queried using SPARQL,<sup>2</sup> the standard query language for the Semantic Web. The model used in the next sections to represent and query trajectory data is based on the concept of property graphs.

Several data models to perform analytical queries on graphs have been proposed. GraphOLAP (Chen *et al.* 2009), conceptually, is a framework for online analytical processing (OLAP) on a set of homogeneous graphs, based on splitting the graph into a

---

<sup>1</sup><https://www.w3.org/RDF/>

<sup>2</sup><https://www.w3.org/TR/rdf-sparql-query/>

collection of snapshots that are aggregated in two ways, called Informational and Topological OLAP aggregations. GraphCube (Zhao *et al.* 2011) provides a framework for computation and analysis on OLAP cubes using the different levels of aggregation of a graph. Gómez *et al.* (Gómez *et al.* 2019) use graph databases to represent semantic trajectory data based on places of interest (PoIs), that is, a collection of trajectories represented as routes between context-defined PoIs rather than actual geographic points (Parent *et al.* 2013).

### 2.3. Graph and Relational Databases

The comparison between relational databases and graph databases has been studied to a limited extent, given that graph database technology is relatively novel. Vicknair *et al.* (Vicknair *et al.* 2010) compare MySQL against Neo4j through a simple database schema and relatively simple queries. A similar study was carried out by Batra and Tyagi (Batra and Tyagi 2012), also using MySQL and early versions of Neo4j. Both studies, however, discuss very simple queries. Regarding spatiotemporal data, Makris *et al.* (Makris *et al.* 2019) compare MongoDB, a document NoSQL database against PostgreSQL, not only for querying but also for data preparation tasks. Gómez *et al.* (Gómez *et al.* 2019) compare graph and relational databases for storing and querying trajectory data, concluding that in most of the queries, the former perform better because they take advantage of the native data storage, in particular for path traversal. The latter is the only study that compares both models for queries that can exploit the natural representation of the model at hand. The present paper works along the same lines, since the river system representation is naturally a network, which can benefit from the native graph data storage of Neo4j in particular, and graph databases, in general. The study is presented in Section 4.

## 3. Data Acquisition and Pre-processing

This section details the data sources used in the paper, and the pre-processing work carried out in order to get the data ready to be exploited. The process includes several non-trivial steps that are worth discussing. First, the data sources are described. Then, the process that transforms the data into a graph containing the river system information is studied in detail.

### 3.1. Data sources

The Flemish environmental agency (whose acronym in Dutch is VMM), produces the “Vlaamse Hydrografische Atlas” (VHA), a data set comprised of shapefiles containing all the rivers in Flanders, and the watersheds the rivers are part of. This data set does not contain ponds and other water bodies. The VHA is maintained by the VMM, and new versions are released every three months. The data set contains geometric data where the rivers in Flanders are represented as line segments, and includes the flow direction of each segment. The main attributes in the data set are (the names of the properties are in Dutch):

- **Vhas**, a unique number that each river segment gets assigned by VMM. This number can be seen and used as an ID of the segment.

- **Catc**, the category to which the segment belongs. All rivers are divided into categories, which range from 0 up to 9, with 0 representing the biggest waterways and 9 the smallest ones.
- **Lengte**, the (precomputed) length of the segment.
- **Geom**, the geometry of the river segment. Most of the the time, it is a multi-line (polygonal) geometry.
- **Naam**, the name of the river
- **Strmgeb**, the name of catchment area
- **Beknaam**, the name of an administrative sub area of the catchment. It can be seen as a broad drainage area.
- **Lblkwal**, the intended quality of the water in the segment, for example “drinkable water”.

The VHA data set includes more properties of the segments, not included here for the sake of space, but included in the databases that are created for this work. All properties and their description can be found the documentation supplied along with the shape-files. Additionally, the OpenStreetMap information is used, since it is considered here as correct and up to date, in general, for Belgium.<sup>1</sup>

Specifically, for the tests reported here, the VHA data set from 7 August 2020 is used.<sup>2</sup>

### 3.2. *Preparing the data set*

The VHA described in the previous section must be processed to produce data that can be used for analysis and prediction. This process is comprised of two steps: (a) Create the relationships between river segments and (b) Fix the errors that may have occurred.

#### 3.2.1. *Creating relationships between river segments*

The representation of the overall water flow must be added to the data set, since the data contains the flow direction within each segment, but not over multiple segments. A new relation is defined encoding this overall flow information. The terminology of the segments needs to be defined first. When water is flowing from one segment to the next one, the two segments involved are called **source** and **target**, respectively. The former is the segment where the water is coming from; the target segment is the one where the water is flowing to. In other words, it can be said that the the target segment follows the source segment for downstream flows. Now that the naming for the two segments involved is known, a relation **flows-to(A,B)**, can be defined as a binary relation where A and B are the IDs of the segments. The relation consists of all tuples  $(a, b)$  where  $a$  is a source segment ID and  $b$  is a corresponding target segment ID.

In order to create this relation, each segment has to be matched with all the other segments, to check whether or not the water flows directly from one segment to the other. The main idea is that the endpoint of the line geometry of the source segment is taken, if there is a starting point of another segment’s line geometry that matches the endpoint, the second segment is a target segment for the source one. These pair of segments can then be added to the **flows-to** relation. This is done for all segments in the VHA, after which the **flows-to** relation represents the complete system flow.

---

<sup>1</sup><https://openstreetmap.be/en/>

<sup>2</sup><http://www.geopunt.be/catalogus/datasetfolder/020a452d-8cd2-41b7-9c64-2be367668837>

It is worth noting that not every segment has a follow-up segment<sup>3</sup>. For example, there are segments that end up in the sea, or just stop in some special cases. This does not always mean that a river stops at the end of that segment; the river can, for example, cross the Flemish border and subsequently not have any follow-up segment in the data set. Also, segments do not always have exactly one follow-up segment, since a river can split into two or more rivers that all flow on and, possibly, join again. In this situation, the endpoint of the segments will fall together with more than one starting point. Therefore, the **flows-to** relation can contain multiple tuples for a specific segment and this should be taken into account when devising algorithms for the search of flow paths, and also for the creation of the database itself.

We remark that, in general, the **flows-to**-graph of a river system is acyclic, since naturally flowing water cannot flow from one location via some path to that same location (if the river system includes pumps, this might be different). Therefore, a **flows-to**-graph is a directed acyclic graph (or DAG).

It has been mentioned that the VHA data set is delivered as a shapefile where all segments and their properties are stored. In order to add the flow information, the file is loaded into a spatial relational database, namely PostgreSQL, equipped with the PostGIS extension<sup>1</sup>. This table is denoted **wlas**. From it, the **flows-to** table is created as:

```
1 CREATE TABLE
2 flows_to(source_segment bigint, target_segment bigint);
```

The new table can be filled using:

```
1 INSERT INTO flows_to(source_segment, target_segment)
2 SELECT a.vhas, b.vhas
3 FROM vlas a, vlas b
4 WHERE ST_StartPoint(b.geom) = ST_EndPoint(a.geom);
```

This query cannot be directly executed over the VHA data set after it is imported into the Postgres database. The reason for that is that the geometries in the VHA shapefile, and thus in the database, are multi-line geometries and the **ST\_StartPoint()** or the **ST\_EndPoint()** functions cannot take a multi-geometry as input. Therefore, the multi-line geometries must be converted to a single line geometry. The following statements create a new column **geomS** in the table **wlas**, with type line geometry defined using the map projection 31370 (which is the “Belgian Lambert 72” projection), and then convert each multi-line segment into a single line segment. After this pre-processing of the VHA data, the query above can be executed. We note that the usage of “b.geom” and “a.geom” needs to be replaced with the name of the newly created column, in this example “b.geomS” and “a.geomS”.

```
1 ALTER TABLE vlas
2 ADD COLUMN geomS geometry(LineString, 31370);
3 UPDATE vlas SET geomS = ST_LineMerge(geom);
```

### 3.2.2. Fixing errors

Some errors encountered during the creation of the data set need to be fixed, to obtain a usable database. These are discussed next.

---

<sup>3</sup>For a given segment, source, the corresponding targets are considered to be follow-up segments.

<sup>1</sup><https://postgis.net>



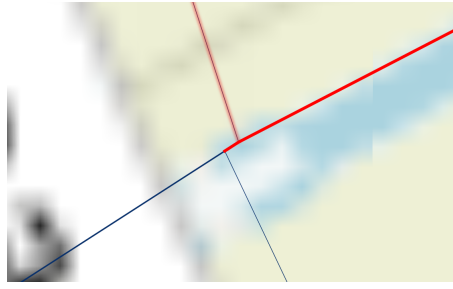


Figure 2. An example of a group of segments where the endpoint does not match exactly with other segments start or end points.

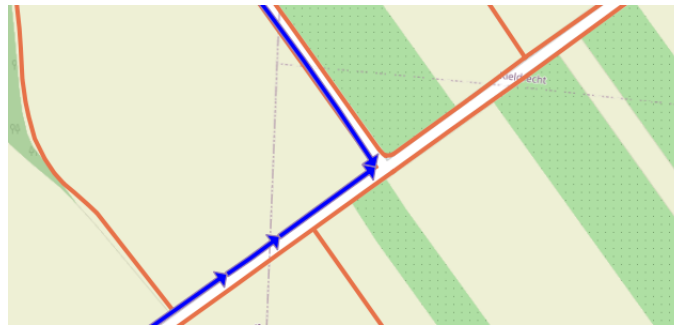


Figure 3. An example of two segments colliding in the VHA caused by and incorrect direction. The blue arrows indicate the direction encoded in the VHA.

**Segments that do not match** Up until now two segments were defined as source and target segments if their ending and starting points coincide. However, if the difference between the two points is relatively small, they may still represent the same physical location (see Figure 2). To overcome such small mismatches, the comparison of the points needs to be relaxed, allowing a *tolerance* for considering two points to be the same. This can be addressed as follows for the original **flows-to** relation:

```

1 WHERE ST_StartPoint(b.geom) = ST_EndPoint(a.geom) OR
2 ST_DWithin(ST_StartPoint(b.geomS),ST_EndPoint(a.geomS),1);

```

Here, the tolerance is set to one meter (the “Belgian Lambert 72” reference system implies meters). The value of the tolerance that is allowed depends on the problem. For the VHA data set, one meter is adopted, based on empirical tests with different values. Using a tolerance also entails that there is a higher possibility of encountering *false positives*, meaning that two segments could appear to be matching although this is not the case. However, with the adopted value this is kept to a reasonable chance.

**Incorrect directions** In the real world, rivers often have a direction associated with them. A direction incorrectly encoded may lead to colliding segments or non-matching starting and ending points, as illustrated in Figure 3.

The solution entails the following steps:

- (1) Execute the standard **flows-to** relation.
- (2) Find all unmatched segments and select the ones that have incorrect-direction issue.
- (3) Store these segments in a temporary table, revert the direction, and add the information back to the created relation.

***Incomplete or unexpected data*** Two further issues must be mentioned:

- Isolated lines or segments; and
- Border crossing and re-entering.

The first one is rather unexpected in a data set that represents rivers or water streams. In the case of the VHA, it can happen that one line geometry or a small set of line geometries do not connect to the rest of the segments in the data set. Those isolated segments form a static water body, like a small pond, where mainly runoff water is captured until it drains into the ground. By definition, the VHA only charts bodies of water or streams that have a flowing character. The segments, therefore, are unexpected data in the VHA data set. The influence of those cases is negligible on the overall data set because they only influence the results if a query, asking for a certain path (discussed in Section 5), starts in one such segment. Furthermore, this occurs with the smallest segments.

Border crossing segments occur because there are administrative boundaries to a region which do not match to the natural boundaries existing in the landscape. The VHA data set contains the rivers, brooks and ditches of the region Flanders. However, the data set, does not include the parts that are outside the region of Flanders. The so-called border segments are the segments that leave the charted area and may re-enter further down the stream. The border segments do not have any downstream segments in the **flows-to** relationship because such a segment does not exist, according to the assumptions made. The problem created by this gap at the border (the water that leaves the border segments, ends up in another segment) is not captured by the present solutions. The distance between those segments can be a few meters but also a few kilometres. In this paper, it is assumed that the water that leaves the region does not re-enter in another stream. This entails that the information of the overall water flow in those cases is lost for good and cannot be captured in the final data set.

## 4. Storing the River Graph in Relational and Graph Databases

In this section, relational and graph databases are analysed with respect to their functional and modelling dimensions. First, since the analysis of rivers involves spatial data, capability to handle these data is discussed first. Then, the representation of the problem using the relational and graph data models is studied. Finally, typical recursive queries over the graph and relational database representations are discussed next along two dimensions: intuitiveness and performance. In this sense, the questions to answer are: how easy and intuitive would be to write typical queries, and how fast will they run? Query performance is studied in Sections 5 and 6.

### 4.1. *Handling spatial and non-spatial data*

The data used in this paper are mainly based on the VHA, which is distributed using shapefiles. Each version of the VHA is a new shapefile including the complete VHA and thus all geometries of the rivers in Flanders. Data are imported into the database and then converted to a suitable representation for querying.

Neither PostgreSQL nor Neo4j are geographical databases, which means that they do not have out-of-the-box support for geometric data like the one in the VHA data

set. However, this support can be added to them through extensions. In the case of PostgreSQL, PostGIS is the extension that adds support for geometric data, through a wide range of geometric functions<sup>1</sup>. PostGIS also provides functions allowing importing a shapefile through the *shp2pgsql* functionality<sup>2</sup>. At the time of writing this work, Neo4j releases do not provide functionalities to import shapefiles. Like in PostgreSQL, the overall flow information must be created after the data are imported (although, in PostgreSQL, PostGIS and pgRouting<sup>3</sup> provide topology creation functions to facilitates this process). There is also a software library that provides interaction with OpenStreetMaps (OSM)<sup>4</sup>, and includes a scalable importer which takes advantage of Neo4j spatial indices, and also provides some functions for routing analysis. In addition, a new spatial library mentioned above, contains algorithms for spatial analysis<sup>5</sup>, although to a much lesser extent, compared with the functionality provided by PostGIS. Finally, the APOC library, that comes with the current Neo4j versions (4.x at the time of writing)<sup>6</sup> contains functions for geocoding over OSM (other map services can be configured). In summary, compared to Postgres, Neo4j so far, lacks a wide range of spatial functionality.

## 4.2. Data model

The typical way of performing routing or path-finding analysis would be to store data in relational databases, over which SQL queries could be run. These queries are aimed at finding paths in the network, aggregating data with respect to some dimensions (e.g., time, river category, and so on), or querying data with respect to some geographical feature, location, or point of interest (PoI). A problem with this approach, particularly with the huge volumes of data available nowadays, is the difference between the way in which data are modelled and stored (this was called “impedance mismatch” above). Given that the river topology can be considered a graph, storing river data using relations may seem unnatural. Especially, since current database technology provides solutions that allow storing graphs in native form, as mentioned in Section 2. Relational and graph data models also come with high-level query languages.

For the problem under study, rivers are modelled as a sequence of segments, connected to each other. This is the typical case of recursive relationships, extensively studied in database conceptual modelling. (Dullea and Song 1999) give a taxonomy of this kind of relationships. The translation of recursive relationships to the relational model is straightforward and also well-studied. Thus, following traditional database theory, the river system is represented as follows. There is a table to store the segments information, such as ID and properties:

$$wlas(vhas, name, \dots).$$

The attribute *vhas* is used as the identifier of the segment, and called ID. There is also a table containing the binary relation *flows-to* is used, as discussed in Section 3.2.1,

---

<sup>1</sup>[https://postgis.net/docs/PostGIS\\_Special\\_Functions\\_Index.html](https://postgis.net/docs/PostGIS_Special_Functions_Index.html)

<sup>2</sup><https://postgis.net/docs/manual-1.4/ch04.html#id419979>

<sup>3</sup><http://pgrouting.org>

<sup>4</sup><https://github.com/neo4j-contrib/osm>

<sup>5</sup><https://github.com/neo4j-contrib/spatial-algorithms>

<sup>6</sup><https://github.com/neo4j-contrib/neo4j-apoc-procedures>

where for each segment the follow-up segments are stored:

$$flowsto(source, target).$$

The `source` and `target` columns contain the IDs of the segments (`vhas`).

In Neo4j, segments are represented as nodes, with label `:Segment` (and their corresponding properties), and the relation between the nodes is called `:flowsTo`, defined as follows: there is a relation `:flowsTo` from node  $A$  to node  $B$  if the water is able to flow to segment  $B$  from segment  $A$ .

We note that in both models, the reverse flow can be addressed when querying, therefore adding the inverse relation, namely `:comesFrom`, is not actually needed to indicate a flow from node  $B$  to node  $A$ .

### 4.3. *Expressing Recursive Queries over Relational and Graph Databases*

Typical queries required by the problem under study are of the form: “Where can the water flow to?” (downstream query) and “Where does the water come from?” (upstream query). Based on these queries, other computations can be performed, like height and speed of the flow, pollution spread models, and many more. We note that these are recursive queries, which are computationally expensive, since they often require computing the transitive closure of the underlying graph, a well-known problem in database theory (For example, see the classic works by (Bancilhon and Ramakrishnan 1986), and (Li and Ross 1993)). Actually, the worst-case time complexity for computing the transitive closure of a directed graph is  $O(n \cdot e)$ , where  $n$  is the number of the nodes, and  $e$  is the number of the edges. The space complexity is  $O(n^2)$ . As an example, a classic algorithm is proposed by (Schmitz 1983). It follows that this is also a hard problem in graph databases. However, this paper shows that the graph representation would better take advantage of the structure of the river system in order to query the database efficiently.

With the layout of the data defined in Section 4.2, the SQL downstream query from a starting segment, with ID `id_startsegment`, can be written as (the upstream query is analogous, and omitted due to space restrictions):

```

1 WITH RECURSIVE outcome(source, target) AS (
2     (SELECT source, target
3      FROM flowsto
4      WHERE source = id_startsegment)
5     UNION
6     SELECT flowsto.source, flowsto.target
7     FROM outcome, flowsto
8     WHERE flowsto.source = outcome.target)
9 SELECT DISTINCT target FROM outcome;
```

Cypher is, like SQL, a high-level, declarative, programming language. It is specifically designed for graph structures, and is the language that comes with the Neo4j database. It uses nodes and relations as first-class citizens, although the output to a query can be a graphs or a set of tuples. The Cypher query that computes the downstream query showed in SQL above, reads:

```

1 MATCH (N:Segment)-[:flowsTo*]->(M:Segment)
2 WHERE N.vhas = id_startsegment
3 RETURN DISTINCT M.vhas;

```

Note that both are recursive queries computing the transitive closure of the graph, and returning the nodes in the graph that can be reached starting from a given one. That is, the queries do not output the paths, but just the reachable segments. In the case of SQL, listing the paths would be even more complex. For example, the query below computes the paths to each reachable segment:

```

1 WITH RECURSIVE outcome(source, target, path) AS (
2     (SELECT flowsto.source, flowsto.target,
3         ARRAY[flowsto.target]
4     FROM flowsto
5     WHERE flowsto.target = id_startsegment)
6 UNION
7     SELECT flowsto.source, flowsto.target,
8         outcome.path || Array[flowsto.target]
9     FROM outcome, flowsto
10    WHERE flowsto.target = outcome.source
11        AND flowsto.target <> All(path))
12 SELECT DISTINCT path FROM outcome;

```

In the case of Cypher, to compute and list the paths, it suffices to write:

```

1 MATCH path= (N:Segment {vhas:id_startsegment})-
2             -[:flowsTo*]->(M:Segment)
3 RETURN DISTINCT path;

```

It can be seen that the structure of the Cypher query is far less complicated and more intuitive than its SQL counterpart, since it takes advantage of the graph structure. In this case, a basic `MATCH .. WHERE .. RETURN` structure suffices to express a recursive query. This is mainly because Cypher is developed as a query language for graphs and recursion is typical in these cases. The `(N:segment)-[:flowsTo*]→(M:segment)` pattern selects all nodes *M* that are reachable by following one or more edges in the graph, traversing the graph using the `:flowsTo` relation. In addition, the APOC library contains many functions that can be used to compute the query above in a more efficient way, using breadth-first and depth-first algorithms for expanding the nodes. An example of such a query is is:

```

1 MATCH (n:Segment {vhas:6033614})
2 CALL apoc.path.expandConfig(n,
3     {relationshipFilter: "flowsTo>", minLevel: 1})
4 YIELD path AS path
5 RETURN path;

```

The `expandConfig` function expands the nodes of a graph, computing all the paths between a node and all the other ones in the graph. Moreover, most of the time, the structure of the river system is a tree (recall that the hydrological models introduced in Section 2, consider a river system as a tree rather than a graph). This allows using functions that compute the (directed) spanning tree of the starting node, which is even more efficient. This function expands a spanning tree reachable from the start node following a relationship up to a certain level adhering to the label filters indicated as

arguments. The nodes returned, collectively, form a spanning tree. This is studied in detail in the next section.

## 5. Querying the River Database

This section discusses a collection of queries over the rivers database. The collection of queries was composed after consultation of several hydrologists. The queries are designed as a starting point for real life challenges as "Where does an observed pollution come from?", "Where will an observed pollution go to? When will it arrive there?" etc. These queries are then run over Neo4j and PostgreSQL, and the results reported in Section 6. The queries are expressed in Cypher and SQL, respectively. However, for the sake of space, only the former are shown here, since SQL is a well-known language, and the work is focused on graph databases. Nevertheless, the SQL queries are included in Appendix 7. For clarity, the queries are organised into classes that account for their main characteristics. To allow an adequate comparison of Cypher and SQL queries, the required output formats are indicated for each query.

**Queries of Type 1 [Aggregation and similarity queries]** The queries in this class are typical *à la* graph OLAP (Gómez *et al.* 2020) kind of aggregate queries. Aggregations are performed over different properties used as categories and metrics. For example, Query 5.1 just uses the segment's length as a metric, while Query 5.2 aggregates this metric by segment category. Query 5.3 takes the length of the segments and compares them against the length of a given node, in order to obtain segments with similar lengths. The output formats are: for Query 5.1 a float number, for Query 5.2, a tuple of the form (key, length), and for Query 5.3, a list of (ID, length) pairs.

Query 5.1 *Compute the average segment length.*

```
1 MATCH (n:Segment)
2 RETURN avg(n.length) AS avglength
```

Query 5.2 *Compute the average segment length by segment category.*

```
1 MATCH (n:Segment)
2 RETURN n.catc as category, avg(n.length)
3 AS avglength order by category asc
```

Query 5.3 *Find all segments that have a length within a 10% margin of the length of segment with ID 6020612.*

```
1 MATCH (n:Segment {vhas:6020612})
2 WITH n.length as length
3 MATCH (m:Segment)
4 WHERE m.length < length*1.1 and m.length > length*0.9
5 RETURN m.vhas, m.length;
```

**Queries of Type 2 [Network Topology]**

This class of queries addresses the computation of metrics of the river network configuration. Although the queries include aggregation (like the ones of Type 1), they are included in this class, according with their main functional meaning.

Query 5.4 *For each segment find the number of incoming and outgoing segments.*

The output of this query is a set of tuples of the form (ID, #in, #out). The query reads in Cypher as follows:

```

1 MATCH (src:Segment)-[:flowsTo]->(n:Segment)-[:flowsTo]
2   ->(target:Segment)
3 RETURN n.vhas as nodenbr, COUNT(DISTINCT src) as segIn,
4        COUNT (DISTINCT target) as segOut

```

Query 5.5 *Find the segments with the maximum number of incoming segments.*

The output of this query is a list of segment IDs and an integer representing the maximum number of incoming segments.

```

1 MATCH (n:Segment)
2 OPTIONAL MATCH (src:Segment)-[:flowsTo]->(n)
3 WITH n, COUNT(distinct src) as indegree
4 WITH COLLECT ([n, indegree]) as tuples,
5      MAX(indegree) as max
6 RETURN [t in tuples WHERE t[1] = max | t]

```

The `OPTIONAL` statement works like a relational outer join. The `COLLECT` statement aggregates the results in a list of pairs, to which list comprehension functions are then applied. The elements in the list, with values equal to the maximum are returned.

Query 5.6 *Find the number of splits in the downstream path of segment 6020612.*

The output of this query is an integer number indicating the number of splits found.

```

1 MATCH (n:Segment {vhas:6020612})
2 CALL apoc.path.spanningTree(n,{relationshipFilter:
3   "flowsTo>", minLevel: 1}) YIELD path AS pp
4 UNWIND NODES(pp) as p
5 MATCH (p)-[:flowsTo]->(r:Segment)
6 WITH p, count(DISTINCT r) as co WHERE co > 1
7 RETURN count(p)

```

Here, the `spanningTree` function from the APOC library is used. This function computes all simple paths that can be reached starting from a node in the graph, using breath-first search by default. This is done visiting nodes only once. The `relationshipfilter` is `"flowsTo>"`, indicating that the path must traverse only this relation, in downstream direction. The function can be parametrised in many ways, for example, indicating the minimum and maximum levels in the path (here, the latter is omitted). A collection of paths is returned (`pp`), which is then flattened as a table with the `UNWIND` statement. All reachable nodes are obtained. For each node in this table, it is tested if this node has more than one outgoing segments. If this is the case, there is a split. The node with `vhas:6020612` is chosen for the test because it is one of the farthest from the sea, thus its flow downstream is one of the longest ones.

Query 5.7 *Find the number of in-flowing segments in the downstream path of segment 6020612.*

The output of the query is an integer giving the number of in-flowing segments found. An in-flowing segment is a segment that ends on the downstream path, but which is not a part of the path itself. That is, a segment that contributes to the flow of a given one.

```

1 MATCH (n:Segment {vhas:6020612})
2 CALL apoc.path.spanningTree(n,{relationshipFilter:
3   "flowsTo>", minLevel: 1}) YIELD path AS pp
4 WITH [p in NODES(pp) | p.vhas] as ids
5 UNWIND ids as id
6 WITH collect(DISTINCT id) as ids
7 MATCH (s:Segment)-[:flowsTo]->(p)
8 WHERE NOT s.vhas in ids AND p.vhas <> 6020612
9       AND p.vhas in ids
10 RETURN count(DISTINCT s) as inflows

```

This query is similar to Query 5.6, also using the `spanningTree` function. List comprehension is used to obtain the node identifiers.

Query 5.8 *Determine if there is a loop in the downstream path of segment 6031518.*

Sometimes, when the level of the sea turns higher than normal, the sea may get into the river flow and reverse its direction. Moreover, anthropogenic influences, such as barriers, dams and sluices, can create loops in the system. From a modelling point of view, in these cases, the graph will contain a cycle. This query finds out if this is the case in the graph under study. This also shows that, in order to get a realistic modelling, the tree representation does not suffice, and a model like the one proposed in this paper is needed. The output of the query is a Boolean.

```

1 MATCH (n:Segment {vhas:6031518})
2 CALL apoc.path.spanningTree(n, {relationshipFilter:
3   "flowsTo>", minLevel: 1}) YIELD path AS pp
4 WITH [p in NODES(pp) | p] as nodelist
5 UNWIND nodelist as p
6 CALL apoc.path.expandConfig(p,
7   {relationshipFilter:"flowsTo>", minLevel: 1,
8   terminatorNodes:[p], whitelistNodes:nodelist})
9 yield path as loop
10 RETURN count(loop) >0 as loops

```

This query needs some explanation, that will also be used later. In this case, not only the `spanningTree` function is used, but also the `expandConfig` function. The left-hand side of Figure 4 shows the representation of the river as segments. Each edge represents a river segment, starting in one node and ending in another one. The representation that was chosen for the graph is depicted on the right-hand side. Here, a segment becomes a node, for example, the segment *c*, running from nodes 3 to 4, becomes the node *c*. It can be seen that segment *g*, for instance, receives flow from two incoming segments, namely *e* and *f*. If, for example, *a* is the starting segment, the `spanningTree` function would only capture one of the paths, the one which is first found by the algorithm. On the other hand, the `expandConfig` function finds all the paths. In a tree representation this problem would not appear, and the second `CALL` would not be needed, greatly simplifying



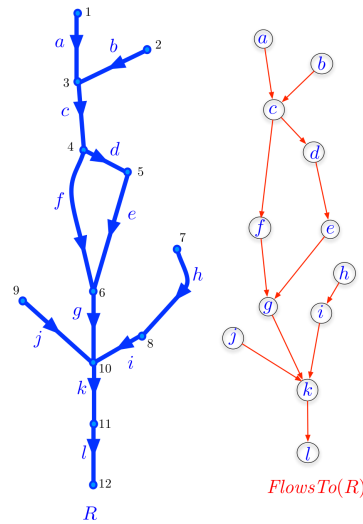


Figure 4. Nodes with more than one incoming flows.

the query. This function has a high computational cost, and should be used only if needed. For example, when the user only needs to obtain the nodes that can be reached from a certain one, `spanningTree` should be used, since it is very efficient.

**Remark 1:** *Discussion on expressiveness.* The Cypher queries above are, in general, simpler than their SQL equivalent (shown in Appendix 7), in particular for Queries 5.5 through 5.8. Writing the SQL code for the latter requires an expert knowledge, while, even though the Cypher equivalents are not trivial, they basically require to know the existence of the right functions. Further, when the `expandConfig` function is not required, the queries turn out to be very simple. In the case of SQL, the queries basically do not change under for the two situations above.

**Queries of Type 3 [Path aggregation]** The queries in this class aggregate a metric along a path. The length of a segment is used here, although for this scenario, the average flow, or of any parameter reported by a sensor, could be used.

**Query 5.9** *Find all paths downstream from the given start segment.*

There is no aggregate function in this query, the aggregation is given by the output, consisting in the IDs of the segments that can be reached from a given one, and the list of IDs of the corresponding paths.

```

1 PROFILE
2 MATCH (n:Segment {vhas:6020612})
3 CALL apoc.path.expandConfig(n, {relationshipFilter:
4     "flowsTo>", minLevel: 1}) YIELD path AS pp
5 UNWIND NODES(pp) as p
6 MATCH (p)-[:flowsTo]->(r:Segment)
7 WITH r, count(DISTINCT p) as co WHERE co > 1
8 WITH collect(r) as pc
9 MATCH (n:Segment {vhas:6020612})
10 CALL apoc.path.expandConfig(n,{relationshipFilter:
11     "flowsTo>", minLevel:1,endNodes:pc}) YIELD path AS pp

```

```

12 WITH [p in NODES(pp) |p.vhas] AS nodelist
13     WHERE size(nodelist) > 0
14 RETURN nodelist[size(nodelist)-1] as id, nodelist
15
16 UNION ALL
17
18 MATCH (n:Segment {vhas:6020612})
19 CALL apoc.path.spanningTree(n,{relationshipFilter:
20     "flowsTo>", minLevel: 1}) YIELD path AS pp
21 UNWIND NODES(pp) as p
22 MATCH (p)-[:flowsTo]->(r:Segment)
23 WITH r, count(DISTINCT p) as co WHERE co = 1
24 WITH collect(r) as pc
25 MATCH (n:Segment {vhas:6020612})
26 CALL apoc.path.spanningTree(n,{relationshipFilter:
27     "flowsTo>", minLevel:1,endNodes:pc}) YIELD path AS pp
28 WITH [p in NODES(pp)|p.vhas] AS nodelist
29 RETURN nodelist[size(nodelist)-1] as id, nodelist;

```

This query requires a trick, to make it possible to run in standard hardware. Since the `expandConfig` function is extremely costly, and the `spanningTree` function is very efficient for reachability, the former is only applied to compute the paths where there is more than one possible path for reaching a segment. This is computed in the upper subquery. The parameter `endNodes:pc` in the functions tell the algorithm to only expand the nodes in this list. The lower subquery uses the `spanningTree` function to compute the paths where there is only one way to reach the segment. The terms `UNION` and `UNION ALL` return the union of the results, without and with duplicates, respectively. This solution would probably not be efficient in a highly interconnected social network, since the `expandConfig` function computes all the paths between a node and all the other ones in the graph, which is computationally very expensive. On the other hand, the `spanningTree` function stops when it finds a path between the node being expanded and each other one. However, it is assumed that river networks are much less interconnected than a typical social network, and therefore it should work well, as shown in the experiments reported in Section 6.

Query 5.10 *Find the branches of downstream flow starting at a given position (identified by a segment's `vhas`), together with the length and number of segments of each branch.*

The output is a collection of tuples of the form: (target segment ID, # of hops, length).

```

1 MATCH (n:Segment {vhas:6020612})
2 CALL apoc.path.spanningTree(n,{relationshipFilter:
3     "flowsTo>", minLevel: 1}) YIELD path AS pp
4 UNWIND NODES(pp) as p
5 MATCH (p)-[:flowsTo]->(r:Segment)
6 WITH r, count(DISTINCT p) as co WHERE co = 1
7 WITH collect(r) as pc
8 MATCH (n:Segment {vhas:6020612})
9 CALL apoc.path.spanningTree(n,{relationshipFilter:
10     "flowsTo>", minLevel: 1,endNodes:pc}) YIELD path AS pp
11 WITH [p in NODES(pp) |p.vhas] AS nodelist,
12 reduce(longi= tofloat(0),n IN nodes(pp)|longi+n.length)

```

```

13         AS segLen,
14     reduce(longi= 1,n IN nodes(pp)| longi + 1) AS nbrSeg
15     RETURN  nodelist[size(nodelist)-1] as id, nbrSeg, segLen;
16
17     UNION
18
19     MATCH (n:Segment {vhas:6020612})
20     CALL apoc.path.spanningTree(n, {relationshipFilter:
21         "flowsTo>", minLevel: 1}) YIELD path AS pp
22     UNWIND NODES(pp) as p
23     MATCH (p)-[:flowsTo]->(r:Segment)
24     WITH r, count(DISTINCT p) as co WHERE co > 1
25     WITH collect(r) as pc
26     MATCH (n:Segment {vhas:6020612})
27     CALL apoc.path.expandConfig(n,{relationshipFilter:
28         "flowsTo>", minLevel:1,endNodes:pc}) YIELD path AS pp
29     WITH [p in NODES(pp) |p.vhas] AS nodelist,
30     reduce(longi = tofloat(0),n IN nodes(pp)|longi+n.lengte)
31         AS segLen,
32     reduce(longi = 1,n IN nodes(pp)| longi + 1) AS nbrSeg
33     RETURN  nodelist[size(nodelist)-1] as id, nbrSeg, segLen;

```

This is similar to the previous query, except for the aggregation of the lengths and number of segments. The `reduce` function computes the value resulting from the application of an expression on each successive element in a list, and accumulates these results as it proceeds. This allows computing the length of each branch (using, in this case, the property `lengte`) and the number of hops.

Query 5.11 *Find the length, the number of segments and the IDs of the segments, of the longest branch of upstream flow starting from a given segment.*

The output is a set of triples of the form (ID, length, # of segments). In this case the length is returned in meters.

```

1  PROFILE
2  MATCH (n:Segment {vhas:6020612})
3  CALL apoc.path.expandConfig(n,{relationshipFilter:
4      "<flowsTo", minLevel: 1}) YIELD path AS pp
5  WITH reduce(longi= tofloat(0), n IN nodes(pp)| longi
6      + tofloat(n.lengte)) AS blength, Length(pp) as
7      alength, [p in NODES(pp) |p.vhas] AS nodelist
8  WITH blength, alength, nodelist[size(nodelist)-1] as id
9  WITH id, max(blength) as ml,
10     collect([id,blength,alength]) as coll
11  WITH id, ml, [p in coll WHERE p[0]= id
12     AND p[1]=ml|p[2]] AS lhops
13  UNWIND lhops as hops
14  RETURN id,ml,hops order by id desc;

```

In this case, the upstream flow is requested. Therefore, the relationship filter now is "`<flowsTo`", indicating that the direction is reversed. This is why there is no need to specify and create the reversed `flows-to`, `comes-from`, relation in the graph. The tricky

part in this query is to solve the cases where the longest physical branch is not the one with the maximum number of hops arriving to the same segment. The function `expandConfig` is used to compute all the alternative paths, and then `reduce` is used to compute the length of each branch. List comprehension is finally used to keep only the tuples that correspond to the branch of maximum length.

Query 5.12 *How many paths exist between two given segments X and Y?*

The output is an integer indicating the number of paths. This case is illustrated by the flow between segments `c` and `g` in Figure 4. To capture this case, again, the function `expandConfig` must be used.

```

1 MATCH (n:Segment {vhas:6020612}),
2     (m:Segment {vhas: 7036554})
3 CALL apoc.path.expandConfig(n,
4     {relationshipFilter:"<flowsTo", minLevel: 1,
5     terminatorNodes:[m]}) yield path as pp
6 RETURN count(pp) as paths

```

**Remark 2:** *Discussion on expressiveness.* Queries in this class are quite complex to write, in both, Cypher and SQL, except Query 5.12, which in Cypher only requires a function call. For the rest of the queries, complexity arises mainly from the situation depicted in Figure 4, which is a very particular case. Otherwise the queries become simpler (although not trivial, of course).

**Queries of Type 4 [Queries with conditions over paths]** These queries only traverse certain branches of the rivers, indicated by conditions over properties of the paths or segments.

Query 5.13 *Find all branches starting at a given segment, reachable traversing the river Scheldt.*

The output is the ID of each final segment, and all the paths that lead to it.

```

1 PROFILE
2 MATCH (n:Segment {vhas:6020612})
3 CALL apoc.path.expandConfig(n,{relationshipFilter:
4     "flowsTo>", minLevel: 1}) YIELD path AS pp
5 UNWIND NODES(pp) as p
6 MATCH (p)-[:flowsTo]->(r:Segment)
7 WITH r, count(DISTINCT p) as co
8 WHERE co > 1
9 WITH collect(r) as pc
10 MATCH (n:Segment {vhas:6020612})
11 CALL apoc.path.expandConfig(n, {relationshipFilter:
12     "flowsTo>", minLevel:1,endNodes:pc}) YIELD path AS pp
13 WITH [p in NODES(pp) WHERE p.strmgeb ="Schelde" |p.vhas]
14     AS nodelist WHERE size(nodelist) > 0
15 RETURN nodelist[size(nodelist)-1] as id, nodelist
16
17 UNION ALL
18

```

```

19 MATCH (n:Segment {vhas:6020612})
20 CALL apoc.path.spanningTree(n,{relationshipFilter:
21     "flowsTo>", minLevel: 1}) YIELD path AS pp
22 UNWIND NODES(pp) as p
23 MATCH (p)-[:flowsTo]->(r:Segment)
24 WITH r, count(DISTINCT p) as co
25 WHERE co = 1
26 WITH collect(r) as pc
27 MATCH (n:Segment {vhas:6020612})
28 CALL apoc.path.spanningTree(n,{relationshipFilter:
29     "flowsTo>",minLevel: 1,endNodes:pc})
30     YIELD path AS pp
31 WITH [p in NODES(pp) WHERE p.strmgeb ="Schelde"|p.vhas]
32     AS nodelist
33 RETURN nodelist[size(nodelist)-1] as id, nodelist;
34
35

```

Since the query asks for all the paths, and not only for the segments, again, the `spanningTree` function is not enough, and `expandConfig` must be used. The statement “[`p in NODES(pp) WHERE p.strmgeb = "Schelde" | p.vhas`]” keeps only the branches of the selected river. Experiments (not reported here, for the sake of space) have proven that this option is more efficient than including a parameter indicating a whitelist of the segments to be traversed.

*Query 5.14 List the length, the number of segments and the IDs of the segments of the the branches starting from a given segment, that are part of the river Scheldt.*

The output are the triples (ID, length, # of segments) for each segment (only the shortest path information). The query is similar to Query 5.13, except for the final part. The computation of the paths is done analogously to the previous query. Thus, for the sake of space only the final part is shown.

```

1 MATCH (n:Segment {vhas:'6020612'})
2 CALL apoc.path.spanningTree(n,{relationshipFilter:
3     "flowsTo>", minLevel: 1}) YIELD path as pp
4     .....
5 WITH [p in NODES(pp) WHERE p.strmgeb ="Schelde" |p.vhas]
6     AS nodelist, reduce(longi= tofloat(0),n IN nodes(pp)|
7     CASE WHEN n.strmgeb ="Schelde" THEN longi + n.lengte
8     ELSE longi END) AS length, reduce(longi= 1,n
9     IN nodes(pp)|CASE WHEN n.strmgeb ="Schelde"
10    THEN longi + 1 ELSE longi END) AS segCount
11 RETURN nodelist[size(nodelist)-1] as id, segCount, length
12

```

The `reduce` statements compute the lengths of the segments and the number of segments in each path. The statement “`CASE WHEN n.strmgeb = "Schelde" THEN longi + tofloat(n.lengte) ELSE longi END`” is used to aggregate only the requested branches in the `reduce` statement. We note that this solution captures all the alternative paths when there is more than one way of reaching a certain node.

**Remark 3:** *Discussion on expressiveness.* Queries in this class are, as it could be seen,

very complex. Query 5.13, in SQL is much simpler, but Query 5.14 requires a deep knowledge of SQL programming.

### *Queries of Type 5* [Spatial queries]

Finally, a class of queries including spatial data are proposed.

Query 5.15 *Find all segments reachable from the segment closest to the Antwerpen Groenplaats<sup>1</sup>.*

The output is a list of segment IDs, no path information is required.

```

1 CALL apoc.spatial.geocodeOnce('Groenplaats
2     Antwerpen Flanders Belgium')
3     YIELD location as ini
4 MATCH (n:Segment)
5 WITH n, ini, distance(
6     point({longitude:n.source_long, latitude:n.source_lat}),
7     point({longitude:ini.longitude, latitude:ini.latitude})
8 ) as d
9 WITH n, d order by d asc limit 1
10 CALL apoc.path.spanningTree(n,
11     {relationshipFilter:"flowsTo>", minLevel: 1})
12     YIELD path as pp
13 UNWIND NODES(pp) as p
14 RETURN p.vhas;
```

Here, the APOC function `geocodeOnce` is used to find the starting point, from which the reachable segments are computed. Antwerpen's Groenplaats is taken as the reference. Then, Cypher's built-in `distance` function computes the distance between Groenplaats and the closest river segment. The rest, is analogous to the previous queries.

Query 5.16 *Find the segments that belong to the downstream path and that are at most at 3 km of the start segment, together with the minimum distance from the start to the segment.*

The output is a list of segment IDs, and the length of the shortest path, in meters. Since the minimum distance is required, again, the `expandConfig` function must be used. Only the portion of the query related with the computation of the distance is shown, the rest is analogous to the previous queries.

```

1 MATCH (n:Segment {vhas:6020612})
2 CALL apoc.path.spanningTree(n, {relationshipFilter:
3     "flowsTo>", minLevel: 1})
4     YIELD path AS pp
5 ...
6 ...
7 CALL apoc.path.expandConfig(n, {relationshipFilter:
8     "flowsTo>", minLevel:1,endNodes:pc}) YIELD path AS pp
9 UNWIND NODES(pp) AS p
10 WITH distance(point({longitude:n.source_long,
```

---

<sup>1</sup>The "Groenplaats" is the main square in the city of Antwerp.

```

11   latitude:n.source_lat})), point({longitude:p.source_long,
12   latitude:p.source_lat})) as dist, p WHERE dist < 3000
13 RETURN DISTINCT p.vhas, min(dist)
14
15 UNION
16 ...
17 ...
18 MATCH (n:Segment {vhas:6020612})
19 CALL apoc.path.spanningTree(n,{relationshipFilter:
20   "flowsTo>", minLevel: 1,endNodes:pc}) YIELD path AS pp
21 UNWIND NODES(pp) AS p
22 WITH distance(point({longitude:n.source_long,
23   latitude: n.source_lat})), point({longitude:
24   p.source_long, latitude: p.source_lat}))
25   as dist, p WHERE dist < 3000
26 RETURN DISTINCT p.vhas, min(dist);

```

In this case, the `distance` function is used to compute which segments are at less than 3 km from the starting point.

**Remark 4:** *Discussion on expressiveness.* Here, comparing the Cypher queries against the SQL and PostGIS queries in Appendix 7, it appears clear that the degree of maturity of spatial capabilities of PostGIS gives SQL a clear edge over the graph alternative. Spatial support is still needed in graph databases.

## 6. Experimental Evaluation

The queries in Section 5 are run over the Neo4j database which is designed and populated as described in Section 4.2. Furthermore, in order to compare performance against the relational alternative, the queries are written in the SQL language, and executed over a PostgreSQL database. For the fairness of the comparison, the type of the output, as well as the results of the SQL queries, are the same as the ones corresponding to the Cypher queries in Section 5. Both databases are fully indexed in order to obtain the best possible query performance. Indices are created over all attributes that are mentioned in the queries (the segment identifiers, `strngeb`, `lengte`, `catc`, etc.). Neo4j provides two classes of indices: Native B-tree and full-text search indices. In this work, native B-tree indices are used. Figure 5 shows the index configuration used for Neo4j.

In PostgreSQL, the tables and indices are stores in the same tablespace. The index type is the default B-tree for all indices. For example, for the `source` attribute in the `wlas` and `flowsto` tables:

```

1 CREATE INDEX edgesour
2   ON public.wlas USING btree
3   (source ASC NULLS LAST)
4   TABLESPACE pg_default;

1 CREATE INDEX flowsto_source_idx
2   ON public.flowsto USING btree
3   (source ASC NULLS LAST)
4   TABLESPACE pg_default;

```

id	name	state	populationPercent	uniqueness	type	entityType	labelsOrTypes	properties	provider
2	"source_idx"	"ONLINE"	100.0	"NONUNIQUE"	"BTREE"	"NODE"	["Segment"]	["source"]	"native-btree-1.0"
7	"source_lat_idx"	"ONLINE"	100.0	"NONUNIQUE"	"BTREE"	"NODE"	["Segment"]	["source_lat"]	"native-btree-1.0"
6	"source_long_idx"	"ONLINE"	100.0	"NONUNIQUE"	"BTREE"	"NODE"	["Segment"]	["source_long"]	"native-btree-1.0"
8	"strmgeb_idx"	"ONLINE"	100.0	"NONUNIQUE"	"BTREE"	"NODE"	["Segment"]	["strmgeb"]	"native-btree-1.0"
3	"target_idx"	"ONLINE"	100.0	"NONUNIQUE"	"BTREE"	"NODE"	["Segment"]	["target"]	"native-btree-1.0"
4	"target_lat_idx"	"ONLINE"	100.0	"NONUNIQUE"	"BTREE"	"NODE"	["Segment"]	["target_lat"]	"native-btree-1.0"
5	"target_long_idx"	"ONLINE"	100.0	"NONUNIQUE"	"BTREE"	"NODE"	["Segment"]	["target_long"]	"native-btree-1.0"

Figure 5. Neo4j index configuration.

Table 1.: Number of nodes and edges in the Rivers graph database.

Type	Name	Size (#)
Node	Segment	61,777
Edge	flowsTo	65,428
<b>Total</b>	<b># Objects</b>	<b>126,205</b>

The starting node reporting in this study is chosen as follows. For downstream flows, the starting segment is chosen close to the start of the flow. For queries analysing upstream flow, starting segments close to the end of the flow are chosen. Although several segments were considered as candidates, only a representative one is reported in this work.

### 6.1. Experiments setup

For the Neo4j database, the number of nodes and edges are given in Table 1. For the PostgreSQL database, the table from where the edges are obtained, called: `wlas`, has 61,777 tuples, and the table `flowsto`, containing overall flow information, 65,428 tuples. The queries are run on a machine with a i7 7700 processor at 2.8GHz, with 32 GB of RAM and a 1 TB disk. The execution times reported are the averages of five runs of each experiment.

### 6.2. Discussion

Table 2 summarises the test results. The last column on the right gives the ratio between the execution times on Neo4j and PostgreSQL. The best execution times for each query have been highlighted in bold font. When the value is set to  $\infty$  this means that the query ran for more than 10 minutes without finishing. The discussion that follows is organised considering the query classes defined in Section 5.

The results show that almost all queries run much faster in Neo4j. Although these results could be expected for transitive-closure like queries, surprisingly, queries of Type 1 (aggregate queries) written in Cypher also outperformed SQL queries, except for Query 5.2. For topological queries (Type 2), Cypher clearly outperforms SQL except for two of the queries. Likewise, performance is, in some cases, orders of magnitude better in Cypher for queries of Types 3 and 4, that means, path queries, which encode the



Table 2.: Execution times for the example queries.

		<b>3</b>	<b>4</b>	
Type of Query	Query	Neo4j (msec)	Postgres (msec)	3 / 4
Aggregation & similarity	5.1	<b>79</b>	94	0.84
Aggregation & similarity	5.2	111	<b>103</b>	1.07
Aggregation & similarity	5.3	<b>96</b>	116	0.83
Network Topology & pattern	5.4	<b>14</b>	258	0.05
Network Topology& pattern	5.5	<b>184</b>	193	0.95
Network Topology& pattern	5.6	<b>35</b>	51	0.69
Network Topology& pattern	5.7	319	<b>47</b>	6.79
Network Topology& pattern	5.8	<b>663</b>	2200	0.30
Path Aggregation & pattern	5.9	<b>1740</b>	$\infty$	N/A
Path Aggregation & pattern	5.10	<b>1820</b>	$\infty$	N/A
Path Aggregation & pattern	5.11	<b>711</b>	47000	0.015
Path Aggregation & pattern	5.12	<b>1</b>	47	0.02
Conditions over paths	5.13	<b>1914</b>	11300	0.17
Conditions over paths	5.14	<b>1596</b>	$\infty$	N/A
Spatial	5.15	<b>388</b>	613	0.55
Spatial	5.16	26038	<b>48</b>	542

computation of the reachability in the graph and conditions and/or aggregations over the paths. Also surprising is the result for queries of Type 5, where spatial functions are used. In this case, however, it is necessary to point out that spatial capabilities for Neo4j are not even close to the ones of PotGIS, as it was already commented earlier. Nevertheless, results are quite good (although of course, far from conclusive). Also in this case, we note that in Query 5.15 the coordinates are computed with the build-in OSM service whereas in PostgreSQL they are hardcoded into the query, and even in this case, performance is better for Neo4j.

Another point that is worth a discussion, is the comparison between using the `spanningTree` function to compute the nodes reachable from a given one, against the simple Cyphers's built-in transitive closure computation (the `*` function). The latter is orders of magnitude worse. However, the `expandConfig` function, which is needed when all the paths must be returned, and not just the nodes reachable from a certain one, is not as efficient, since it computes all the paths in the transitive closure.

Also, the analysis of the queries in Section 5 suggests that, in general, expressing queries in a graph-based high-level language, results in simpler, more concise, and more intuitive expressions than their SQL equivalent. However, there are situations, typical in NoSQL databases, where the way in which a query is written impacts on the performance. This particularly occurs when all the paths must be computed, and the river system is modelled as a graph. When only a segment's reachability is required, or the river system can be modelled as a tree, or alternative paths are not needed, the Cypher expressions can be highly simplified, while SQL queries still require computing the transitive closure of the relation.

## 7. Conclusion and Future Work

This paper uses a real-world case, based on the Flemish river system, to study the plausibility of using graph databases to represent, store, and query river data. The work also presents a traditional relational database implementation, and compares both alternatives. The data preparation tasks are described, as well as the data models used. Finally, a collection of queries are defined, and executed over the PostgreSQL and Neo4j

databases, expressed in SQL and in Cypher, the high-level query languages for both databases, respectively. The queries are run and the results discussed and reported.

The study suggests that river systems, and other kinds of transportation networks, can be modelled as a graph, and implemented using graph databases, over which queries are, in general, more easily expressed using high-level graph query languages, in this particular case, Cypher. The results also show that queries involving path computation run overall faster over graph databases, since their underlying data structures are designed to achieve fast path traversal. Opposite to this, a relational representation requires writing recursive queries to compute the transitive closure of the graph, which impacts on query efficiency, since the relational representation does not capture the graph structure appropriately, a problem known in database modelling as “impedance mismatch”. Five types of queries were studied, including aggregate, path, and spatial queries. Only three out of sixteen queries delivered better performance in the relational version. In particular, in path computation, where the graph representation is crucial, the difference reaches orders of magnitude in favour of Cypher. However, it is worth noting that these results were obtained through the algorithms provided in Neo4j libraries, not with Cypher’s built-in transitive closure computation. Nevertheless, long path traversals like the ones required in this problem, are clearly not appropriately handled by the relational model, since they require multiple self-joins of the table containing the relationships between the river segments. It must be mentioned that intensive, advanced SQL query tuning was not the scope of this work. Rather it was the intention to investigate the feasibility of using graph databases to model river networks. In summary, the results obtained in this work suggest that graph databases can become a good alternative for analysing large volumes of river data, like the ones in the IoW project.

Future work is mainly oriented at scaling this problem for larger volumes, for which parallel processing may be needed. There are many parallel processing graph databases (e.g., GraphFrames<sup>1</sup>, Janusgraph<sup>2</sup>) that may take advantage of the characteristics of graphs like the ones studied here. Even Neo4j has recently presented a scalable version in the cloud. Other future work consists in a generalization to other transportation networks like road networks, computer networks, sewage networks, heat networks, among other ones.

## Acknowledgements

Erik Bollen was supported by the *Bijzonder Onderzoeksfonds* (BOF) from UHasselt with reference BOF20OWB27 and by VITO with project reference 2010478.

Alejandro Vaisman was partially supported by Project PICT 2017-1054, from the Argentinian Scientific Agency.

---

<sup>1</sup>[https://graphframes.github.io/graphframes/docs/\\_site/index.html](https://graphframes.github.io/graphframes/docs/_site/index.html)

<sup>2</sup><http://janusgraph.org/>

## References

- Ahani, A., Shourian, M., and Rahimi Rad, P., 2018. Performance Assessment of the Linear, Nonlinear and Nonparametric Data Driven Models in River Flow Forecasting. *Water Resources Management*, 32 (2), 383–399.
- Angles, R., 2012. A Comparison of Current Graph Database Models. In: *Proceedings of ICDE Workshops*, Arlington, VA, USA, 171–177.
- Angles, R., *et al.*, 2017. Foundations of Modern Query Languages for Graph Databases. *ACM Comput. Surv.*, 50 (5), 68:1–68:40.
- Angles, R., 2018. The Property Graph Database Model. In: *Proceedings of the 12th Alberto Mendelzon International Workshop on Foundations of Data Management, Cali, Colombia, May 21-25, 2018*.
- Angles, R. and Gutierrez, C., 2008. Survey of graph database models. *ACM Comput. Surv.*, 40 (1), 1:1–1:39.
- Bancilhon, F. and Ramakrishnan, R., 1986. An Amateur’s Introduction to Recursive Query Processing Strategies. In: C. Zaniolo, ed. *Proceedings of the 1986 ACM SIGMOD International Conference on Management of Data, Washington, DC, USA, May 28-30, 1986* ACM Press, 16–52.
- Batra, S. and Tyagi, C., 2012. Comparative Analysis of Relational And Graph Databases. *International Journal of Soft Computing and Engineering*, 2.
- Brouwers, J., *et al.*, 2015. *MIRA Climate Report 2015*. Technical report, VMM, Aalst.
- Chen, C., *et al.*, 2009. Graph OLAP: a multi-dimensional framework for graph data analysis. *Knowl. Inf. Syst.*, 21 (1), 41–63.
- da Silva, M.C.T., *et al.*, 2015. SWOT: A Conceptual Data Warehouse Model for Semantic Trajectories. In: *Proceedings of the ACM Eighteenth International Workshop on Data Warehousing and OLAP, DOLAP 2015, Melbourne, VIC, Australia, October 19-23, 2015*, 11–14.
- Daltio, J. and Medeiros, C.B., 2015. HYDROGRAPH: EXPLORING GEOGRAPHIC DATA IN GRAPH DATABASES. *Revista Brasileira de Cartografia*, 68 (6).
- Demir, I. and Szczepanek, R., 2017. Optimization of river network representation data models for web-based systems. *Earth and Space Science*, 4 (6), 336–347.
- Dullea, J. and Song, I., 1999. A Taxonomy of Recursive Relationships and Their Structural Validity in ER Modeling. In: *Conceptual Modeling - ER ’99, 18th International Conference on Conceptual Modeling, Paris, France, November, 15-18, 1999, Proceedings*, Vol. 1728 of *Lecture Notes in Computer Science* Springer, 384–398.
- Fileto, R., *et al.*, 2015. The Baquara<sup>2</sup> knowledge-based framework for semantic enrichment and analysis of movement data. *Data Knowl. Eng.*, 98, 104–122.
- Gobin, A., 2012. Impact of heat and drought stress on arable crop production in Belgium. *Natural Hazards and Earth System Science*, 12 (6), 1911–1922.
- Gómez, L.I., Kuijpers, B., and Vaisman, A.A., 2019. Analytical queries on semantic trajectories using graph databases. *Trans. GIS*, 23 (5), 1078–1101.
- Gómez, L.I., Kuijpers, B., and Vaisman, A.A., 2020. Online analytical processing on graph data. *Intell. Data Anal.*, 24 (3), 515–541.
- Hartig, O., 2014. Reconciliation of RDF\* and Property Graphs. *CoRR*, abs/1409.3288.
- Heuvelmans, G., Muys, B., and Feyen, J., 2004. *Analysis of the spatial variation in the parameters of the SWAT model with application in Flanders, Northern Belgium*. Technical report 5.
- Li, Z. and Ross, K., 1993. *On the Cost of Transitive Closures in Relational Databases*. Technical report CUCS-004-93, Columbia University.

- Makris, A., *et al.*, 2019. Database system comparison based on spatiotemporal functionality. In: B.C. Desai, D. Anagnostopoulos, Y. Manolopoulos and M. Nikolaidou, eds. *Proceedings of the 23rd International Database Applications & Engineering Symposium, IDEAS 2019, Athens, Greece, June 10-12, 2019* ACM, 21:1–21:7.
- Parent, C., *et al.*, 2013. Semantic trajectories modeling and analysis. *ACM Comput. Surv.*, 45 (4), 42:1–42:32.
- Robinson, I., Webber, J., and Eifré, E., 2013. *Graph Databases*. O'Reilly Media.
- Ruback, L., *et al.*, 2016. Enriching Mobility Data with Linked Open Data. In: *Proceedings of the 20th International Database Engineering and Applications Symposium, IDEAS 2016, Montreal, QC, Canada, July 11-13, 2016*, 173–182.
- Schmitz, L., 1983. An improved transitive closure algorithm. *Computing*, 30 (4), 359–371.
- Solomatine, D.P. and Ostfeld, A., 2008. Data-driven modelling: Some past experiences and new approaches. *Journal of Hydroinformatics*, 10 (1), 3–22.
- Vicknair, C., *et al.*, 2010. *A Comparison of a Graph Database and a Relational Database A Data Provenance Perspective*. .
- Zhao, P., *et al.*, 2011. Graph cube: on warehousing and OLAP multidimensional networks. In: *Proceedings of the 2011 ACM SIGMOD International Conference on Management of Data*, Athens, Greece ACM, 853–864.

## Appendix A

### SQL Queries

Query 5.1 *Compute the average segment length.*

```
1 SELECT avg(lengte) AS avglength
2 FROM edges;
```

Query 5.2 *Compute the average segment length by segment category.*

```
1 SELECT catc AS category, avg(lengte) AS avglength
2 FROM wlas
3 GROUP BY catc;
```

Query 5.3 *Find all segments that have a length within a 10% margin of the length of segment with ID 6020612.*

```
1 SELECT vhas, lengte
2 FROM wlas
3 WHERE lengte <= 1.1*
4     (SELECT lengte from edges WHERE vhas=6020612)
5     AND lengte >= 0.9*(SELECT lengte
6     FROM edges
7     WHERE vhas=6020612)
8
```

Query 5.4 *For each segment find the number of incoming and outgoing segments.*

```
1 SELECT segments.vhas, count(DISTINCT flowstoB.source) AS segIn,
2     count(DISTINCT flowstoA.target) AS segOut
3 FROM wlas as segments, flowsto as
```

```

4         flowstoA, flowsto as flowstoB
5 WHERE segments.vhas = flowstoA.source
6         AND segments.vhas = flowstoB.target
7 GROUP BY segments.vhas;

```

Query 5.5 *Find the segments with the maximum number of incoming segments.*

```

1 SELECT target as vhas, segIn FROM
2     (SELECT target, count(flowsto.source) AS segIn
3     FROM flowsto
4     GROUP BY flowsto.target) AS myTable
5 WHERE segIn = (SELECT max(segIn) FROM
6     (SELECT flowsto.target,
7     count(flowsto.source) AS segIn
8     FROM flowsto
9     GROUP BY flowsto.target) AS tt);

```

Query 5.6 *Find the number of splits in the downstream path of segment 6020612.*

```

1 SELECT count(source) FROM
2 (WITH RECURSIVE outcome(source, target) AS (
3     (SELECT source, target
4     FROM flowsto
5     WHERE source = 6020612)
6     UNION
7     SELECT flowsto.source, flowsto.target
8     FROM outcome, flowsto
9     WHERE flowsto.source = outcome.target )
10 SELECT source, count(target) AS segOut
11 FROM outcome
12 GROUP BY source) AS myTable
13 WHERE segOut > 1;

```

Query 5.7 *Find the number of in-flowing segments in the downstream path of segment 6020612.*

```

1 SELECT sum(diff)
2 FROM (
3     SELECT myTable.target, count(source)-segIn as diff
4     FROM
5         (WITH RECURSIVE outcome(source, target) AS (
6             (SELECT source, target
7             FROM flowsto
8             WHERE source = 6020612)
9             UNION
10            SELECT flowsto.source, flowsto.target
11            FROM outcome, flowsto
12            WHERE flowsto.source = outcome.target )
13         SELECT target, count(source) AS segIn
14         FROM outcome
15         GROUP BY target) AS myTable, flowsto
16     WHERE myTable.target = flowsto.target
17     GROUP BY myTable.target, segIn) AS secTable;

```

Query 5.8 *Determine if there is a loop in the downstream path of segment 6031518.*

```

1 WITH RECURSIVE outcome(source, target, again, path) AS (
2   (SELECT source, target, 0, ARRAY[source]
3     FROM flowsto
4     WHERE source = 6031518)
5   UNION
6   SELECT flowsto.source, flowsto.target,
7     CASE WHEN flowsto.source <> All(path) THEN 0
8       ELSE 1 END, outcome.path || Array[flowsto.source]
9   FROM outcome, flowsto
10  WHERE flowsto.source = outcome.target AND
11    outcome.again <> 1)
12  SELECT count(source)>0 FROM outcome where again=1;

```

Query 5.9 *Find all paths downstream from the given start segment.*

```

1 WITH RECURSIVE outcome(source, target, path) AS (
2   (SELECT flowsto.source, flowsto.target,
3     ARRAY[flowsto.source]
4   FROM flowsto
5   WHERE flowsto.source = 6020612)
6   UNION
7   SELECT flowsto.source, flowsto.target, outcome.path
8     || Array[flowsto.source]
9   FROM outcome, flowsto, was
10  WHERE flowsto.source = outcome.target AND
11    flowsto.source <> All(path))
12  SELECT json_agg(array_to_json(outcome.path)) AS paths
13  FROM outcome
14  WHERE 0=(SELECT count(target)
15    FROM outcome as cin
16    WHERE outcome.target=cin.source)
17  GROUP BY outcome.target;

```

Query 5.10 *Find the branches of downstream flow starting at a given position (identified by a segment's vhas ID), together with the length and number of segments of each branch.*

```

1 WITH RECURSIVE outcome(source, target, path, length, segCount)
2 AS (
3   SELECT flowsto.source, flowsto.target, ARRAY[flowsto.source],
4     w1.lengte + w2.lengte, 1
5   FROM flowsto, was as w1, was as w2
6   WHERE flowsto.source = 6020612 and flowsto.source = w1.vhas
7     and flowsto.target = w2.vhas
8   UNION
9   SELECT flowsto.source, flowsto.target, outcome.path
10     || Array[flowsto.source], outcome.length + was.lengte,
11     outcome.segCount + 1
12  FROM outcome, flowsto, was
13  WHERE flowsto.source = outcome.target AND
14    was.vhas = flowsto.target AND flowsto.source <> All(path))

```

```

15  SELECT target, path, length, segCount FROM outcome
16  WHERE 0=(SELECT count(target) FROM outcome as cin
17           WHERE outcome.target=cin.source);

```

Query 5.11 *Find the length, the number of segments and the IDs of the segments, of the longest branch of upstream flow starting from a given segment.*

```

1  WITH RECURSIVE outcome(source, target, path,
2                        length, segCount) AS (
3    (SELECT flowsto.source, flowsto.target,
4           ARRAY[flowsto.target], w1.lengte
5           + w2.lengte, 1
6    FROM flowsto, wlas as w1, wlas as w2
7    WHERE flowsto.target = 6020612 AND
8           flowsto.source = w1.vhas AND
9           flowsto.target = w2.vhas)
10   UNION
11   SELECT flowsto.source, flowsto.target,
12          outcome.path || Array[flowsto.target],
13          outcome.length + wlas.lengte,
14          outcome.segCount + 1
15   FROM outcome, flowsto, wlas
16   WHERE flowsto.target = outcome.source
17          AND wlas.vhas = flowsto.source
18          AND flowsto.target <> All(path))
19  SELECT source, min(length), min(segCount)
20  FROM outcome
21  GROUP BY outcome.source;

```

Query 5.12 *How many paths are there between two given segments X and Y?*

```

1  WITH RECURSIVE outcome(source, target, path) AS (
2    (SELECT flowsto.source, flowsto.target,
3           ARRAY[flowsto.target]
4    FROM flowsto
5    WHERE flowsto.target = 6020612)
6    UNION
7    SELECT flowsto.source, flowsto.target, outcome.path
8           || Array[flowsto.target]
9    FROM outcome, flowsto
10   WHERE flowsto.target = outcome.source AND
11          flowsto.target <> All(path) AND 7036554 <> All(path))
12  SELECT count(DISTINCT path)
13  FROM outcome
14  WHERE 7036554 = Any(path);

```

Query 5.13 *Find all branches starting at a given segment, reachable traversing the river Scheldt.*

```

1  WITH RECURSIVE outcome(source, target, path) AS (
2    (SELECT flowsto.source, flowsto.target,
3           ARRAY[flowsto.source]
4    FROM flowsto

```

```

5 WHERE flowsto.source = 6020612)
6 UNION
7 SELECT flowsto.source, flowsto.target, outcome.path
8     || Array[flowsto.source]
9 FROM outcome, flowsto, was
10 WHERE flowsto.source = outcome.target AND flowsto.source
11     <> All(path) AND was.vhas = flowsto.source
12     AND strmgrb = 'Schelde')
13 SELECT outcome.target, json_agg(array_to_json(path))
14 FROM outcome
15 GROUP BY outcome.target;

```

Query 5.14 *List the length, the number of segments and the IDs of the segments of the branches starting from a given segment, that are part of the river Scheldt.*

```

1 WITH RECURSIVE outcome(source, target, path,
2     length, segCount) AS (
3     (SELECT flowsto.source, flowsto.target,
4         ARRAY[flowsto.source], w1.length + w2.length, 1
5     FROM flowsto, was as w1, was as w2
6     WHERE flowsto.source = 6020612 and flowsto.source =
7         w1.vhas and flowsto.target = w2.vhas)
8     UNION
9     SELECT flowsto.source, flowsto.target,
10         outcome.path || Array[flowsto.source],
11         outcome.length + was.length, outcome.segCount+1
12 FROM outcome, flowsto, was
13 WHERE flowsto.source = outcome.target AND
14     flowsto.source <> All(path) AND was.vhas =
15     flowsto.target AND strmgrb = 'Schelde')
16 SELECT DISTINCT outA.target, outA.length, outB.segCount
17 FROM outcome as outA, outcome as outB
18 WHERE outA.target = outB.target AND
19     outA.length=(SELECT min(length)
20 FROM outcome as c2 WHERE c2.target=outA.target)
21     AND outB.segCount=(SELECT min(segCount)
22 FROM outcome as c3
23 WHERE c3.target=outB.target);

```

Query 5.15 *Find all segments reachable from the segment closest to the Antwerpen Groenplaats.*

```

1 WITH RECURSIVE outcome(vhas) AS (
2     (SELECT was.vhas
3     FROM was
4     ORDER BY ST_Distance(ST_Point(source_long, source_lat),
5         ST_Point(4.4016, 51.2192)) LIMIT 1)
6     -- 51.2192, 4.4016 are coordinates of Groenplaats Antwerpen
7     UNION
8     SELECT flowsto.target
9     FROM outcome, flowsto
10    WHERE outcome.vhas = flowsto.source)

```



```
11 SELECT DISTINCT vhas FROM outcome;
```

Query 5.16 *Find the segments that belong to the downstream path and that are at most at 3 km of the start segment, together with the minimum distance from the start to the segment.*

```
1  WITH RECURSIVE outcome(vhas, path, dist, geom) AS (
2      (SELECT wlas.vhas, ARRAY[vhas], 0.0::double precision, geom
3      FROM   wlas
4      WHERE  vhas = 6020612)
5      UNION ALL
6      SELECT flowsto.target, outcome.path || Array[flowsto.target],
7      ST_Distance(ST_StartPoint(ST_LineMerge(wlas.geom)),
8      ST_StartPoint(ST_LineMerge(outcome.geom))), outcome.geom
9      FROM outcome, flowsto, wlas
10     WHERE outcome.vhas = flowsto.source AND
11           flowsto.target = wlas.vhas
12     AND flowsto.target <> All(path) AND
13           ST_Distance(ST_StartPoint(ST_LineMerge(wlas.geom)),
14           ST_StartPoint(ST_LineMerge(outcome.geom))) < 3000)
15 SELECT vhas, dist FROM outcome;
```