

Network Data Model and BerlinMOD Benchmark

Simone Jandt

Last Update: December 18, 2009

Abstract

In the past several data models for the representation of spatial and spatio-temporal data objects have been developed. Among others we can categorise the data models into data models for objects moving freely in space and data models for moving objects that are constrained by a given network. For both categories several different data models have been developed. Two of them, one data model for objects moving freely in the 2D plane and one for network constrained moving objects have been implemented in SECONDO DBMS. Both of them handling with the history data of moving objects. In this paper we present the comparison of the capabilities of the both data models using the BerlinMOD Benchmark. The results show the advantages of specialised data models for network constrained objects. Therefore we propose an extension of the BerlinMOD Benchmark with an additional query set BerlinMOD/Net including specialised challenges for network data models. This extension enables us to compare the capabilities of different network data models with the BerlinMOD Benchmark.

1 Introduction

In the past several data models for the representation of spatial and spatio-temporal data objects have been developed. Among others we can categorise them into data models for objects moving freely in space and data models for objects which movement is constrained by a network. For both categories several different data models have been presented like [9, 1] for spatio-temporal data objects moving freely in space and [10, 16, 20] for spatio-temporal data objects constrained by networks to name just a few.

For our experiments we choosed one example data model as representative for each of the both categories to compare the capabilities of the both different data models. The category of data models for spatio-temporal objects moving freely in space is represented by the data model presented in [9]. And the category of network data models is represented by the network data model presented in [10]. The advantage of this choice is, that both data models are implemented in the SECONDO DBMS and use the same temporal representation. So we can exclude that the use of different DBMS or temporal representations impact our results.

We used the recently introduced BerlinMOD Benchmark [3] to compare the capabilities of the both data models. The BerlinMOD Benchmark is best to our knowledge the first benchmark for complete spatio-temporal database systems and it is available in the SECONDO DBMS. Furthermore the BerlinMOD Benchmark data model is the data model of freely movement we use for our comparison. So we only have to transform the BerlinMOD Benchmark spatial and spatio-temporal data types once in our network data model representation. This avoids sources of errors and relieves the control for correctness of the query results.

The results of our experiments show that it makes sense to use specialised network data models for network constrained objects, because the network data model outperforms the data model of

free movement in the space in nearly all test cases. The total run time over all queries for the network data model is one-third of the total run time for the data model of free movement in the space. So the further development of network data models seems to be useful. But network data models have their own challenges like network distance computing therefore we propose an extension of the BerlinMOD Benchmark with an additional query set BerlinMOD/NET. This new query set should cover the specialised challenges of network data models. And enable us to compare the capabilities of different network data models also for the specialised network challenges with the BerlinMOD Benchmark.

The rest of the paper is organised as follows. In section 2 we will first give short reminders of the underlying SECONDO DBMS 2.1, the BerlinMOD Benchmark 2.2 and the data models (2.3 and 2.4) used in this paper. In section 3 we present the setup for our experiments in 3.1 and describe our transfer of the BerlinMOD Benchmark in the network data model representation in 3.2 before we show our experimental results in section 3.3. In section 4 we present our proposed extension of the BerlinMOD Benchmark query sets. We conclude the paper with a summation of our results and aspects of future work in section 5.

2 Related Work

As mentioned before in the past other data models for free movement in the space [9, 1] and for network constrained movement have been presented [10, 16, 20]. As well there are some more benchmarks [13, 19] and database systems for spatial and spatio-temporal datatypes [14, 11]. But they don't provide a combination of different implemented and supported data types together with an existing benchmark like the actual SECONDO DBMS. It is self-evident for us to use the SECONDO DBMS in combination with the provided data types and benchmark to compare the capabilities of the both different representations of spatial and spatio-temporal data objects for network constrained objects.

In the next subsections we give short reminders of the used database system (2.1), the used benchmark (2.2), and the used data models (2.3, 2.4). More detailed information for each of them can be taken from the original papers.

2.1 Secondo

The extensible SECONDO DBMS presented in [5, 11] provides a platform for implementing various kinds of data models defining a clean interface between the data model independent system frame and the content of the single data models. The data model independent components and mechanisms are separated from the data model dependent parts. Hence SECONDO can be easily extended by implementing SECONDO algebra modules to provide new data types and operations by the user. The user also may define additional viewers for the graphical user interface or write optimisation rules or cost functions to extend the optimiser. SECONDO is free available in the web [8] and comes with a number of already implemented spatial and spatio-temporal data types and operations including the spatio-temporal data model of free movement in the space 2.3 and the network data model 2.4 (since version 2.8.4). Furthermore the BerlinMOD Benchmark described in 2.2 has been developed in the SECONDO DBMS. For our experiments we used the SECONDO version 2.9.1 with some minor bugfixes in the NetworkAlgebra and the TemporalNetAlgebra.

2.2 BerlinMOD Benchmark

The BerlinMOD Benchmark was recently presented in [3] and is implemented in the SECONDO DBMS. It is available in the web [7] and provides a well defined data-set and queries for the exper-

imental evaluation of different moving object data representations. The BerlinMOD Benchmark emphasises the development of complete systems and simplifies experimental repeatability pointing out the weakness and the potency of the benchmarked systems.

The data-sets of the BerlinMOD Benchmark are created using the street map of Berlin [15] and statistical data about the regions of Berlin [17, 18] as input relations. The created moving objects represent cars driving in the streets of Berlin. This makes it possible to use the data set of the BerlinMOD Benchmark also for network constrained data models. Every moving object has a home node and a work node and every weekday there will be a trip from the home node to the work node in the morning and trip from the work node back to the home node in the afternoon. In the evening and at the weekend randomly chosen cars spend additional trips (one in the evening and up to six at the weekend) to different randomly chosen targets. The number of observed cars and the duration of the observation period can be influenced by the user of the BerlinMOD Benchmark by setting the *scalefactor* to different values in the data generation script. For example at *scalefactor* 1.0 the data generator will create 2000 moving objects observed for 28 days. Each of them sending a GPS-signal every 2 seconds. This simulated signals are simplified so that time intervals when a car doesn't move or moves in the same direction at the same speed are merged into one unit. E.g. if the car holds 8 hours in front of the work node there will be only one entry in the cars history of movement with a time interval of 8 hours instead of 14.400 entries one for each GPS interval.

The BerlinMOD Benchmark provides two different approaches to store the histories of moving objects. On the one hand the object-based approach (OBA) and on the other hand the trip based approach (TBA). In the OBA the complete history for each object is kept together. In the TBA exist several trips for each object. In the OBA there is only one relation (*dataScar*) containing one tuple for each object consisting of the spatio-temporal data of the object (journey), the licence, the type, and the model of the object. In the TBA we have two relations. One of them (*dataMcar*) contains the static data for each object like licence, type, and model together with an object identifier. The other relation (*dataMtrip*) contains for each object identifier several tuples each of them containing a single trip of the moving object, whereby breaks between the trips (e.g. the time the car holds in front of the office) are stored as single (stationary) trips.

Besides the moving objects the BerlinMOD Benchmark provides sets of *QueryPoints*, *QueryRegions*, *QueryInstants*, *QueryPeriods*, and *QueryLicences*, each of them containing 100 objects (points, regions, time instants, time intervals, and licences) used in the benchmark queries.

For this data objects the BerlinMOD Benchmark provides two sets of queries. One set addresses range queries (BerlinMOD/R) and the other one nearest neighbour queries (BerlinMOD/NN). In this paper we focus on the range queries which are the main aspect of the BerlinMOD Benchmark up to now. The query set BerlinMOD/R includes 17 queries selected of the set of possible combinations of the 5 aspects object identity (known / unknown), dimension (standard / spatial / temporal / spatio-temporal), query interval (point / range / unbounded), condition type (single object / object relations), and aggregation (with or without aggregation). We will present the 17 queries in more detail in section 3.2.4 together with our network data model translations.

2.3 BerlinMOD Data Model

The data model used by the BerlinMOD Benchmark is the same data model of freely moving in 2D space presented in [6, 9, 12]. All spatial positions are given in x,y-coordinates. A single spatial position is represented by the data type *point*. A *point* consists of a pair of *real* values interpreted as x,y-coordinates in the 2D plane. Streets are represented by *line* values. A *line* value consists of a set of half segments representing the geometry of the line in the 2D plane. Each half segment consists of two *point* values which are interpreted as start and end point of the half segment. Regions are

represented by the data type *region*. A *region* consists of a set of half segments interpreted as outer (and inner) border of the region in the 2D plane.

All this spatial data types and many standard data types can be lifted to become time dependent *moving* values. For all data types *a* the constructor *moving* creates a new data type *moving(a)* (short form *ma*). A car may be represented by a *mpoint*. A *mpoint* is a *point* changing its position within time. Therefore a *mpoint* consist of a set of units called *unit(point)* (short form *upoint*). Each *upoint* consists of a time interval and two *point* values. The first point represents the position of the *mpoint* at the start of the time interval and the second point represents the position of the *mpoint* at the end of the time interval. The *point* is assumed to move on the straight line between this two points with constant speed. The speed is given by the ratio from the distance of the two points and the length of the time interval of the unit. All units of a *mpoint* must have disjoint time intervals, because a car cannot be at two different positions at the same time. The units are sorted by ascending time intervals. This spatio-temporal data model of *moving* allows us to compute the position of a *mpoint* at every time instant within its definition time. We can also compute the time instant the point passed a given position assumed the *mpoint* ever passes this position. The position of a *point* at a given time instant is represented by a *intime(point)* (short form *ipoint*). A *ipoint* consists of a time instant and a *point* value.

Some other data types of SECONDO which are used in the BerlinMOD Benchmark are shown in table 1.

Data Type	Description
<i>bool</i>	Usual boolean data type.
<i>int</i>	Usual integer number.
<i>real</i>	Usual real number.
<i>instant</i>	A point in time.
<i>periods</i>	A set of disjoint and not connected time intervals.
<i>mbool</i>	A time dependent boolean value. The value within each <i>ubool</i> will be constant <i>TRUE</i> or <i>FALSE</i>
<i>mreal</i>	Time dependent real number. Each unit will be defined by a function of time representing the <i>real</i> value at each time instant.

Table 1: Other Data Types of BerlinMOD Benchmark

2.4 Network Data Model

The central idea of the network data model presented in [10] is that every movement is constrained by a network and every position is given related to this network. Therefore the data type *network* is the central data type in the network data model. All other data types of the network data model are related to a *network* by the unique network identifier that is part of each *network* object. The *network* object contains all spatial information of the represented network in three main relations. The first relation contains the attributes of the routes (streets) like id, route curve, route length, and two boolean flags indicating if the route starts at the lexicographic smaller end point and if the lanes of the route are separated like on German Highways or not. The second relation contains all attributes of the junctions (street crossings) like the identifiers of the first and second route crossing in the junction, the distance of the junction from the start of the first respectively second route, tuple identifiers of the both routes in the routes relation, tuple identifiers of the sections connected by this junction in the sections relation, and a connectivity code telling us which lanes of the two routes are connected by the junction. The third main relation of a *network* object is the sections relation containing the attributes of the sections (street parts between two junctions or a junction and the end of the street) like the route identifier of the route the section belongs to, the tuple identifier of this route in the routes relation, start and end position of the section on the route, section curve, and two boolean flags *startssmaller* and *dual* with the same meaning as in the

routes relation. We introduce four B-Tree indexes and one R-Tree index to support faster query execution. The four B-Trees indicate the route identifier attributes in the three main relations. And the R-Tree indicates the curve attribute of the routes relation. Furthermore there are two sets which provide a fast access from one section to their adjacent sections in driving direction¹.

Single positions in the network are given as gpoint values. Besides the network identifier a gpoint consists of a route identifier, a distance from the start of the route to the position of the gpoint and a side value (*up*, *down*, *none*). The side value is basically necessary if the route is dual, it tells us if the position is reachable from the *up* or the *down* side of the route. For simple routes or positions which are reachable from both sides of the route the side value is always *none*.

Parts of the network, regardless if they represent paths or regions, are given as gline values. Besides the network identifier a gline consists of a set of route intervals, and two boolean flags telling us if the gline is defined and if the set of route intervals is sorted. Each route interval consists of a route identifier identifying the route the route interval belongs to, and the start and the end position from the route interval on this route. In the original paper the route interval includes a additional parameter side like the gpoint does. But this parameter is not part of the implementation yet. We call the set of route intervals sorted if the following conditions are fulfilled:

1. all route intervals are disjoint
2. the route intervals are stored in ascending order of their route identifiers
3. if two disjoint route intervals have the same route identifier the route interval with the smaller start position is stored first
4. for all route intervals in the set the condition: *start position* \leq *end position* holds

We introduced the sorted gline because many algorithms take profit from sorted gline values. For example the computation if a gpoint is inside the gline takes $O(n)$ for unsorted and $O(\log n)$ time for sorted gline values, if n is the number of route intervals in the gline.

Unfortunately not all gline values can be stored sorted. If a gline value represents a path between two gpoint in the network we need the route intervals exactly in the sequence they are used in the path. This will nearly never be a sorted set like defined before. We solved this dilemma by introducing the sorted flag. Every algorithm which can take profit from a sorted gline values checks this flag and uses the corresponding code. We store gline values sorted whenever this is possible to support faster query execution.

Mostly similar to the mpoint of the other data model we implemented a mgpoint. A mgpoint consists of a set of ugpoint with disjoint time intervals. Each ugpoint consists of a time interval and two gpoint values. Every time the mgpoint changes the route or the speed a new ugpoint is written. Each ugpoint is assumed to follow the same route from the start to the end position at the same speed. So accordingly to the mpoint we can compute the network position of the mgpoint at ever time instant within the definition time of the mgpoint as intime(gpoint).

In deviation from the original network data model we extended the implementation of the mgpoint with four additional attributes:

1. The total driven distance
2. A sorted set of route intervals representing the positions ever traversed by the mgpoint
3. A boolean defined flag for the set of route intervals
4. A spatio-temporal minimum bounding box

¹We call two sections adjacent if their lanes are connected by a junction.

The sorted set of route intervals was introduced, because it makes it much faster to decide if a mgpoint ever passed a given place or not. Instead of a linear check of all m ugpoints of a mgpoint we perform a binary scan on the much lower number r of route intervals. This reduces the time complexity from $O(m)$ to $O(\log r)$ for all **passes** operations. Logical this should be done by a sorted gline value but the SECONDO DBMS restricts us to use a sorted set of route intervals instead.

The spatio-temporal minimum bounding box was introduced as parameter to the mgpoint because the computation of this value is very expensive in the network data model. Although each unit of a mgpoint stays on the same route at same speed the motion may follow different spatial directions, e.g. a route may lead uphill in serpentine. Not all this positions must be enclosed by a bounding box computed just from the spatial position of the *start* and *end* position of the unit. Therefore we have always to examine the spatial dimensions of the complete part of the route passed within a unit to compute the units bounding box. All spatial information of the route curve is hidden in the network object. We have to call the route curve from the network object to compute the spatial dimensions of the unit bounding box. If r is the number of routes of the network and h the number of half segments of the traversed part of the route curve passed in a unit we need $O(\log r + h)$ time to compute the bounding box for a single unit. The bounding box of the mgpoint is the union of the bounding boxes of its m units. So the computation of a mgpoint bounding box takes $O(m(h + \log r))$ time. This is very expensive therefore the bounding box of a mgpoint is only computed on demand or if we can get it for free. E.g. we can copy the bounding box of a mpoint if we translate it into a mgpoint without computational effort. But we don't maintain this attribute at every change of the mgpoint. If the mgpoint changes we set the bounding box attribute to be undefined and compute it again on demand if necessary.

3 BerlinMOD and Network Data Model

In the next subsections we first present our experimental setup (3.1) and the transfer from the BerlinMOD Benchmark data-sets and queries in our network data model representation including a description of the indexes we build to support faster query execution (3.2) before we conclude the section with the results of our experiments in 3.3.

3.1 Experimental Setup

For our experiments we used a Standard PC with a Intel Pentium 4 2.93 GHz CPU, 1 GB main memory, 200 GB HDD, and Linux openSUSE 11.1 as operating system. We installed SECONDO version 2.9.1+ and the BerlinMOD Benchmark from the web.

We used the data generating script of the BerlinMOD Benchmark with *scalefactor* value 0.05, 0.2, and 1.0 to generate three data sets with different amounts of data in three different database directories. After that we build the index structures used by the BerlinMOD Benchmark with the script "BerlinMOD_CreateObjects.SEC" delivered with the BerlinMOD Benchmark data for each database and started the benchmark queries for the object oriented and the trip based approach of the BerlinMOD Benchmark on this databases. We saved the results for each database and measured the run times of the queries several times to be sure that the run times measured are free from other influences.

Next we translated the three databases into network data model representation. Therefore we build a network object from the data of *streets relation* and translated all spatial and spatio-temporal data types of the BerlinMOD Benchmark data sets relative to this network object. We also defined some indexes to support faster query execution on the network data model representation and formulated executable SECONDO queries for each BerlinMOD/R query on this network data set. We give a detailed description of the translation steps, indexes and queries in section 3.2.

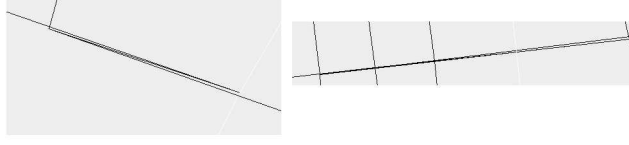


Figure 1: Example Failures in Street Map

After that we started a first run of our network queries and compared the results of our network queries with the results of the BerlinMOD Benchmark queries to ensure that all results are correct. We found some isolated mismatches in some query results, which are caused by the fact that singular route curves in the source data of the street map was not well defined see figure 1 at two places. We corrected the source file “streets.data” delivered by BerlinMOD Benchmark at this two routes and restarted the building of the databases and our experiments from the start. After all results from the network queries corresponded to the results of the BerlinMOD Benchmark results we started the network queries several times to measure the run times analogous to the run times of the benchmark queries. The results of the run time measurement are shown in 3.3.

3.2 Translation of BerlinMOD into Network Data Model

In section 3.2.1 we describe the construction of the central *network* object from the *streets* value of the BerlinMOD Benchmark. This *network* object is necessary to transfer the spatial and spatio-temporal data objects into the network data model representation like described in section 3.2.2. In section 3.2.3 we describe the indexes we build on the data sets in the network data model to support a faster execution of the queries described in section 3.2.4.

A SECONDO script for the network data creation and the translation of the BerlinMOD Benchmark data sets into network data representation and SECONDO scripts with the executable SECONDO BerlinMOD Benchmark queries on the network data model translation can be taken from our website [7].

3.2.1 Create Network Object

We extract the routes data from the *streets* of the BerlinMOD Benchmark data set to create a *network* object for our network data model representation of the BerlinMOD Benchmark. The extracted routes data *r* is used to compute the crossings of the routes of Berlin *j*. In this step the connectivity code for each crossing is set to the maximum value because the data source lacks on information about the connectivity of the crossings. We use *r* and *j* as input relations for the creation of our *network* object *net* representing the streets of Berlin in the network data model.

The network creation algorithm first copies all tuples of *r* to the *routes relation* of *net* and creates the B-Tree index of the route identifiers and the R-Tree of the route curves of the *routes relation* of *net*. Then it copies all tuples of *j* to the *junctions relation* of *net* and adds the *tuple identifiers* for the both routes connected by this junction to the junctions entry in the *junctions relation*. After that two B-Trees indexes for the route identifiers of the first respectively second route in the *junctions relation* are build. Next for every route of the *routes relation* all junctions *jr* on this route are taken to compute for each of this junctions the up and down sections on the route. The up and down sections are inserted into the *sections relation* of *net* and the *tuple identifiers* of the sections are added to the entry of the according junctions in the *junctions relation*. After that the B-Tree index for the route identifiers in the *sections relation* is created and the adjacency lists

of *net* are filled with the adjacent section pairs defined by the *junctions relation*.

If $|r|$ is the number of routes and $|j|$ is the number of junctions. The algorithm needs $O(|r| \log |r|)$ time to copy r to the *routes relation* of *net* and create the trees indexing the *routes relation* of the *net*. The creation of the *junctions relation* and the build of the B-Trees indicating the route identifiers in this relation takes $O(|j| \log |j|)$ time. $O(|r||j|)$ time is needed to fill the *sections relation* and $O(|j|)$ time to fill the *adjacency lists* of *net*. Altogether the complete algorithm needs: $O(|r| \log |r| + |j| \log |j| + |r||j|) = O(|r||j|)$ time to create the *net* from the two input relations r and j .

3.2.2 Translate Spatial and Spatio-Temporal Data

We translate the spatial and spatio-temporal data types of the BerlinMOD Benchmark data set into network objects related to the *network* object *net* described in 3.2.1. The input for all algorithms is a spatial respectively spatio-temporal object and the *network* object to which they should be related. If a input data object is not constrained by the given *network* object the result of the translation is undefined for all translation algorithms described in the following paragraphs.

We start explanation of our translation algorithms with the **point2gpoint** operation, because this operation is used by the other translation algorithms. The algorithm translates a *point* value p into a corresponding *gpoint* value by using the R-Tree index of the *network routes relation* to select the route closest to p and computes the position of p on this route. The *side* value of the resulting *gpoint* is always set to *none*. If r is the number of routes in the *network routes relation* and k is the number of possible candidate routes the worst case complexity of the algorithm is $O(k + \log r)$.

This should be all to translate the *point* values of the *QueryPoints* relation of the BerlinMOD Benchmark into network query positions. But we have a problem with the network data model representation of junctions. In the network data model contrary to the data model of freely moving in space junctions have more than one *gpoint* representation, because they are related to two or more routes. Hence if a junction position is given related to route a we won't detect the junction as passed if a *mgpoint* object passes the junction on route b in all cases, because the definition of **passes** in the network data model is slightly different from the **passes** operation in the BerlinMOD Benchmark data model. Unfortunately all query points of the BerlinMOD Benchmark are junctions. As work around we added a operator **polygpoints**, which returns for every input *gpoint* value gp a stream of *gpoint* values. If gp represents a junction we return all *gpoint* values representing the same junction in the *network* object, otherwise we return only gp in the stream. So we got 221 query *gpoint* values in *QueryPointsNet* for the 100 query *point* values in *QueryPoints* and 22 *gpoint* values in *QueryPoints1Net* for the 10 *point* values of *QueryPoints1* of the BerlinMOD Benchmark. Although we have always to compute the results for the doubled number of query points our network data model outperforms the data model of free movement in nearly all cases by orders of magnitudes.

The second operation **mpoint2mgpoint** translates a *mpoint* value s into a *mgpoint* value t . The main idea of the algorithm is to use the continuous movement of s to reduce computation time. We initialise the algorithm by reading the first unit of s and use the **point2gpoint** operation to find a route in the network containing the *start point* and the *end point* of this unit. We initialise the first unit of t with the computed network values. Then we read each unit of s and try to find the *end point* of the unit on the same route the last unit of s was found. If the *end point* is found on this route we check the direction and speed of the unit. If they are equal to the last unit we extend the actual unit of t to enclose the value of the actual unit of s . If the speed or the moving direction changes we write the actual unit to t and initialise a new unit for t with the network values of the actual unit from s . If the *end point* cannot be found on the same route than the last

unit from s we write the actual unit of t and start a search on the route curves of the adjacent sections to find the route curve that contains the *start point* and the *end point* of the actual unit of s . We initialise a new unit for t with the estimated network values for the actual unit of s and continue with the next unit of s . At least we add the actual network unit to t . The time complexity to find the start values for the first unit is $O(\text{point2gpoint})$. For the next m units of the *mpoint* the worst case complexity is $O(a)$ for each unit if a is the maximum number of adjacent sections. So we get a worst case time complexity of $O(O(\text{point2gpoint}) + ma)$ for the translation of a *mpoint* into a *mgpoint*.

The translation of the *region* values in the *QueryRegions* relation of the BerlinMOD Benchmark into *gline* values of our network data model is done in several steps. First of all we build a single big *line* object from all our network streets. Then we compute for each *region* of the *QueryRegions* the intersection with this big *line* object. At least we translate the resulting *line* objects of the intersection, each representing one *region* of the *QueryRegions* relation, into sorted *gline* values using the *line2gline* operation. The algorithm of the *line2gline* operation takes each *half segment* of a *line* value and computes a corresponding network *route interval* by searching a common route curve for the *start point* and the *end point* of the *half segment* using the *point2gpoint* operation. The computed *route intervals* are sorted, merged and compressed before the resulting *gline* value is returned. If the number of *half segments* of a *line* value is h and the number of resulting compressed *route intervals* is r we get a time complexity of $O(hO(\text{point2gpoint}) + h \log r + r)$ for the whole algorithm. Whereby the summand $h \log r + r$ is caused by the compressing and sorting of the resulting *gline* but as mentioned before in 2.4 we think this time is well invested.

3.2.3 Create Indexes on Network Data Model

After translating all the BerlinMOD Benchmark data sources we are able to create indexes on our network data representation of the BerlinMOD Benchmark data. First we create B-Trees for the *licences* and *moid* attributes of the relations *dataSNcar*, and *dataMNtrip* similar to the indexes created in the BerlinMOD Benchmark for *dataSCcar*, and *dataMCtrip*, because the relations *dataSNcar* and *dataMNtrip* contain the network data model representation of the *dataScar* and *dataMtrip* relation of the BerlinMOD Benchmark. Then we create R-Tree indexes over the spatio-temporal bounding boxes of the *mgpoint* attributes in the *dataMNtrip* and the *dataSNcar* relation. At least we create some specialised network indexes indicating network positions and network-temporal positions of moving objects. Therefore we introduced two and three dimensional *netboxes*. A *netbox* is a degenerated two or three dimensional rectangle. The coordinates of the rectangle are defined to be $x_1 = x_2 = \text{route identifier as real value}$ (The equality of x_1 and x_2 makes the degeneration.), $y_1 = \min(\text{start position, end position})$, $y_2 = \max(\text{start position, end position})$, and, in the three dimensional case, $z_1 = \text{start time as real value}$ and $z_2 = \text{end time as real value}$. For every unit of each *mgpoint* we build a three dimensional *netbox* and for every *route interval* of every *mgpoint* a two dimensional *netbox*. This *netboxes* are used to create R-Trees over the network and network-temporal positions of the *mgpoints* in the network data representation of the BerlinMOD Benchmark.

3.2.4 Translate Benchmark Queries

We developed executable SECONDO queries for each of the 17 BerlinMOD/R queries for the object and the trip based approach using our network indexes to support faster query execution. We had to do this manually because the SECONDO optimiser is not able to optimise SQL-queries on network data model objects yet. In our experiments we tried many different query formulations for each query to get optimal queries delivering the correct result in a minimum of time. The limited space does not allow us to show all our executable SECONDO network queries in detail. As

mentioned before the complete `SECONDO` scripts can be taken from our website. In the following we describe only the algorithms of a view queries in detail.

Every time we need a licence in the result or have a query licence number we have a additional step in the TBA. Because we have to join the *trip* attribute from *dataMNtrip* with the *licence* attribute from *dataMNcar* using the *moid* attribute and the corresponding B-Tree indexes. This will not be repeated at every single query description.

Query 1 and 2 work only on standard attributes. They are formulated analogous to the original queries of the BerlinMOD Benchmark only the relation names and the btrees are changed to match the network data model.

Query 3 selects uses the licence B-Tree to select the ten cars with licences from *QueryLicences1* from *dataSNcar* then the positions of this cars are computed for each of the ten time instants from *QueryInstants1*.

In Query 4 we produce a *netbox* for each of the *QueryPointsNet* and use our specialised netbox R-Tree of the *route intervals* of the *mgpoint* to select the passing vehicles.

In the queries 5, 6, and 10 a retransmission of network objects into spatial respectively spatio-temporal objects of the BerlinMOD Benchmark data model is done. This is caused by the fact that the BerlinMOD Benchmark deals with Euclidean Distances. Euclidean Distances are not very useful in network environments because all objects are restricted to use network paths. Therefore in networks normally the Network Distance is computed. To make the results comparable we retransmit the intermediate results of our network data model into spatial and spatio-temporal data types and use the existing spatial and spatio-temporal Euclidean Distance Functions of the BerlinMOD Benchmark data model for the distance computation in the queries 5, 6, and 10.

Query 5 selects the cars with licences from *QueryLicence1* respectively *QueryLicences2* using the B-Tree over the *Licence* attribute of *dataSNcar* and creates a *line* value from the list of *route intervals* passed by every car. Then the Euclidean Distance between this *line* values is computed for each pair of licences one from *QueryLicences1* and one from *QueryLicences2*. In the TBA we need a aggregation step building the union of the several *mgpoint* belonging to each candidate car. This is done with the *route intervals* because it takes much less time to build the union of the *route intervals* than of the *half segments* in the *line* values representing the same network part.

Query 6 uses the **filter** operation to select the "trucks" from *dataSNcar* (respectively *dataMcar* in TBA) relation. Then the spatio-temporal bounding box of each trip is computed and the spatial dimensions of this box are extended by 5m in every spatial direction. After that the *mgpoint* values are retranslated into *mpoint* values. In a second step each result of the first step is joined with all other results of the first step if the extended bounding boxes intersect, the licences are different and the *mpoint* values have sometimes a distance lower than 10m. The licence pairs of trucks fulfilling this predicate are returned. In the TBA there might be duplicate licence pairs which we have to remove before we return the result.

The first part of the first step of query 7 is almost equal to the selection of cars passing a query point in query 4. The intermediate result is filtered to remove all "not passenger" cars and for every remaining trip the time the trip reaches first the query position is computed for every query position and every candidate trip. In a second step the resulting time instants are grouped by the *id* of the query positions and the minimum time stamp of each group is computed. This minimum time stamp is for every query position the first time it was reached by a car. In the third and last step the licences of the cars reaching the query positions at this first time instant are computed by a join of the results of the first two steps by query position id and the equality of the time stamps.

In query 8 we just select the candidate cars with the licence B-Tree and compute for every car the length of the trip at the query periods in the OBA. In the TBA we have to aggregate over all the distances driven in the single trips by a car within a query period.

For query 9 we compute the length of every trip in every query period, and select the maximum driven distance for every period. In the TBA again we have to do a aggregation of the distances driven from the same car in the same period.

For query 10 we first retranslate every *mgpoint* value of *dataSNCar* into a *mpoint* value. In the OBA we extend the spatial bounding box of each of this both trips by 1.5 m in every direction. Second we select the ten candidate trips given by *QueryLicences1*, retranslate them and extend their spatial bounding boxes. Than we use **symmjoin** to join all trips from the first and the second step where the extended bounding boxes intersect. We filter the pairs that have different licences and are sometimes nearer than 3 m. For this pairs we compute the position of the *mgpoint* at the times the distance between the both *mpoint* has been smaller than 3 m. We return the licence pairs and the positions when they have been closer than 3 m to each other. Again we have to do in the TBA a aggregation of the resulting trips for each licence pair into one. We tried out several indexes to support faster query execution of query 10 including the MON-Tree [4]. But at least this simple form shows the best performance of all.

In query 11 we build a network-temporal query box from the product of *QueryInstant* and *QueryPoints1Net* relation. And use the network-temporal index on *dataSNcar* (respectively *dataM-Ntrip*) to select the resulting trips.

The first step of query 12 is identical with query 11. In a second step a product of the result of the first step with itself is computed and checked for vehicles which have been at the same query point at the same query time instant.

Query 13 first computes a network-temporal box for the query values and then uses this computed boxes to select the candidates with help of the network-temporal R-Tree. In TBA again the resulting *moids* must be mapped to the licences of the cars to generate the result.

Query 14 and 15 work almost similar to query 13 only the query value building the netboxes are different.

Query 16 selects the candidate trips using the licence B-Tree and filters them by the **passes** operation and restricts them to the times they were inside the query region. Then this reduced trips are filtered to be **present** within the query periods and are restricted to the times of the query periods. This is done twice. One time for *QueryLicences1* and one time for *QueryLicences2*. The both results are joined to get the trips of different cars which where at the same period in the same region without meeting each other there in a third step. Again in the TBA we have to do a additional selection from trips with the *moids* belonging to the cars selected before by the licences and to remove duplicates of licence pairs in the same period.

Query 17 again uses the methods from query 4 to find the trips passing a given query point. The passing cars are grouped by the passed query points and the number of cars per query point is computed. In a second step the point with the maximum number of hits is selected and his id and the number of passing cars is returned. In the TBA we have to remove the duplicate cars from the result before computing the hits.

3.3 Benchmark Results

We made several runs for each data amount and each query to get correct average execution times for each query in both data models. The tables 2 and 3 show the resulting run times for each query in seconds. Figure 2 visualises the run time comparison between the different data models and approaches.

	Scalefactor 0.05		Scalefactor 0.2		Scalefactor 1.0	
Query	OBA	TBA	OBA	TBA	OBA	TBA
1	0.173	0.213	0.160	0.152	x.xxx	x.xxx
2	0.006	0.006	0.011	0.007	x.xxx	x.xxx
3	0.751	0.531	1.491	0.557	x.xxx	x.xxx
4	22.365	39.873	158.999	214.883	x.xxx	x.xxx
5	2.351	3.937	5.892	7.607	x.xxx	x.xxx
6	49.187	36.910	246.216	320.830	x.xxx	x.xxx
7	34.860	37.560	514.199	239.386	x.xxx	x.xxx
8	0.893	0.955	2.768	2.583	x.xxx	x.xxx
9	245.831	610.944	741.497	3282.183	x.xxx	x.xxx
10	326.769	141.512	1604.588	1141.317	x.xxx	x.xxx
11	0.186	0.249	1.572	0.818	x.xxx	x.xxx
12	1.843	0.287	118.927	0.349	x.xxx	x.xxx
13	36.509	36.148	350.815	164.260	x.xxx	x.xxx
14	1.372	1.701	14.816	8.656	x.xxx	x.xxx
15	3.143	3.387	115.370	48.872	x.xxx	x.xxx
16	89.351	38.791	62.673	70.790	x.xxx	x.xxx
17	3.600	17.825	179.519	195.433	x.xxx	x.xxx
Total	819.190	970.737	4119.612	5698.686	x.xxx	x.xxx

Table 2: Query Run Times BerlinMOD Benchmark

	Scalefactor 0.05		Scalefactor 0.2		Scalefactor 1.0	
Query	OBA	TBA	OBA	TBA	OBA	TBA
1	0.128	0.161	0.361	0.179	x.xxx	x.xxx
2	0.020	0.008	0.089	0.009	x.xxx	x.xxx
3	0.130	0.158	0.287	0.800	x.xxx	x.xxx
4	0.579	2.853	1.949	16.634	x.xxx	x.xxx
5	1.669	1.785	3.853	4.163	x.xxx	x.xxx
6	13.754	7.955	79.812	57.158	x.xxx	x.xxx
7	6.047	7.801	199.326	195.778	x.xxx	x.xxx
8	0.323	0.318	0.700	1.761	x.xxx	x.xxx
9	58.509	67.456	142.271	168.100	x.xxx	x.xxx
10	157.318	186.050	816.874	1160.634	x.xxx	x.xxx
11	0.233	0.493	6.268	7.996	x.xxx	x.xxx
12	0.313	0.317	0.581	0.577	x.xxx	x.xxx
13	3.019	5.895	62.987	119.839	x.xxx	x.xxx
14	2.559	2.429	22.776	17.628	x.xxx	x.xxx
15	0.632	1.101	27.000	43.437	x.xxx	x.xxx
16	0.905	2.353	1.546	3.621	x.xxx	x.xxx
17	0.294	1.255	2.623	12.884	x.xxx	x.xxx
Total	246.431	288.388	1369.305	1811.198	x.xxx	x.xxx

Table 3: Query Run Times Network Data Model

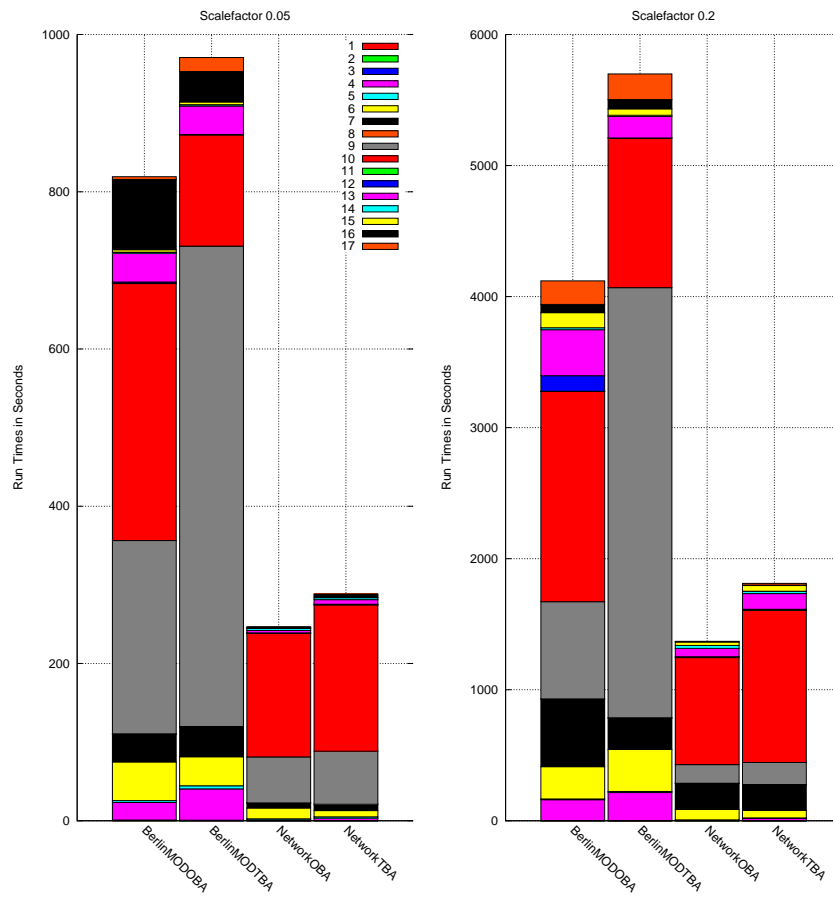


Figure 2: Compared Run Times for Each Query

4 New BerlinMOD Queries

5 Summary and Future Work

Our experiments show that the network data model outperforms the BerlinMOD Benchmark data model in the most cases. The good results of the network data model encouraged us to extend the BerlinMOD Benchmark with a set of queries that enables us to compare the capabilities of different spatio-temporal network data models with respect to the specialised network challenges of this data models.

Actually we have a students practice group implementing a other network data model in SECONDO DBMS. So that we will be enabled to compare the capabilities of the different network data models with the extended BerlinMOD Benchmark in the near future.

References

- [1] Cindy Xinmin Chen and Carlo Zaniolo. *Sqlst a spatiotemporal model and query language*. In *Conceptual Modeling - ER 2000*, volume 1920/2000, pages 111–182. Springer Berlin, Heidelberg, 2000.
- [2] Thomas Behr Christian Düntgen and Ralf Hartmut Güting. *Berlinmod: A benchmark for moving object databases*. Informatik Berichte 340, FernUniversität in Hagen, 2007.
- [3] Thomas Behr Christian Düntgen and Ralf Hartmut Güting. *Berlinmod: A benchmark for moving object databases*. *The VLDB Journal*, 2009.
- [4] Victor Teixeira de Almeida and Ralf Hartmut Güting. *Indexing the trajectories of moving objects in networks*. Informatik Berichte 309, FernUniversität in Hagen, 2004.
- [5] Stefan Dieker and Ralf Hartmut Güting. *Plug and play with query algebras: Secondo-a generic dbms development environment*. In *IDEAS '00: Proceedings of the 2000 International Symposium on Database Engineering & Applications*, pages 380–392, Washington, DC, USA, 2000. IEEE Computer Society.
- [6] Martin Erwig, Ralf Hartmut Güting, Markus Schneider, and Michalis Vazirgiannis. *Spatio-temporal data types: An approach to modeling and querying moving objects in databases*. *Geoinformatica*, 3(3):269–296, 1999.
- [7] Fernuniversität Hagen. *BerlinMOD Benchmark Web Site*, June 2008.
- [8] Fernuniversität Hagen. *Secondo Web Site*, April 2009.
- [9] Luca Forlizzi, Ralf Hartmut Güting, Enrico Nardelli, and Markus Schneider. *A data model and data structures for moving objects databases*. In *SIGMOD '00: Proceedings of the 2000 ACM SIGMOD international conference on Management of data*, pages 319–330, New York, NY, USA, 2000. ACM.
- [10] Hartmut Güting, Victor Teixeira de Almeida, and Zhiming Ding. *Modeling and querying moving objects in networks*. *The VLDB Journal*, 15(2):165–190, 2006.
- [11] Ralf Hartmut Güting, Victor Almeida, Dirk Ansoerge, Thomas Behr, Zhiming Ding, Thomas Hose, Frank Hoffmann, Markus Spiekermann, and Ulrich Telle. *Secondo: An extensible dbms*

- platform for research prototyping and teaching. In *ICDE '05: Proceedings of the 21st International Conference on Data Engineering*, pages 1115–1116, Washington, DC, USA, 2005. IEEE Computer Society.
- [12] Ralf Hartmut Güting, Michael H. Böhlen, Martin Erwig, Christian S. Jensen, Nikos A. Lorentzos, Markus Schneider, and Michalis Vazirgiannis. A foundation for representing and querying moving objects. *ACM Trans. Database Syst.*, 25(1):1–42, 2000.
- [13] Christian S. Jensen, Dalia Tiesyte, and Nerius Tradisauskas. The cost benchmark -comprasion and evaluation of spatio-temporal indexes. In *Database Systems for Advanced Applications*, volume 3882/2006, pages 125–140. Springer Berlin, Heidelberg, 2006.
- [14] Nikos Pelekis, Yannis Theodoridis, Spyros Vosinakis, and Themis Panayiotopoulos. Hermes - a framework for location based data management. In *Advances in Database Technology - EDBT 2006*, volume 3896/2006, pages 1130–1134. Springer Berlin, Heidelberg, 2006.
- [15] S. Rezić. *Berlin Road Map*, 2008.
- [16] Laurynas Speičvcys, Christian S. Jensen, and Augustas Kligys. Computational data modeling for network-constrained moving objects. In *GIS '03: Proceedings of the 11th ACM international symposium on Advances in geographic information systems*, pages 118–125, New York, NY, USA, 2003. ACM.
- [17] Statistisches Landesamt Berlin. *Bevoelkerungsstand in Berlin Ende September 2006 nach Bezirken*, 2008.
- [18] Statistisches Landesamt Berlin. *Interaktiver Stadtatlas Berlin*, 2008.
- [19] Yannis Theodoridis. Ten benchmark database queries for location-base services. *The Computer Journal*, 46(6):713–725, 2003.
- [20] Michalis Vazirgiannis and Ouri Wolfson. A spatiotemporal model and language for moving objects on road networks. In *Advances in Spatial and Temporal Databases*, volume 2121/2001, pages 20–35. Springer Berlin, Heidelberg, 2001.