



DEPARTMENT OF COMPUTER SCIENCE,
FACULTY OF MATHEMATICS, PHYSICS AND INFORMATICS,
COMENIUS UNIVERSITY IN BRATISLAVA

DISTANCE ORACLES FOR TIMETABLE GRAPHS

(Master thesis)

bc. František Hajnovič

Study program: Computer science

Branch of study: 2508 Informatics

Supervisor: doc. RNDr. Rastislav Kráľovič, PhD.

Bratislava 2013



Comenius University in Bratislava
Faculty of Mathematics, Physics and Informatics

THESIS ASSIGNMENT

Name and Surname: Bc. František Hajnovič
Study programme: Computer Science (Single degree study, master II. deg., full time form)
Field of Study: 9.2.1. Computer Science, Informatics
Type of Thesis: Diploma Thesis
Language of Thesis: English
Secondary language: Slovak

Title: Distance oracles for timetable graphs

Aim: The aim of the thesis is to explore the applicability of results about distance oracles to timetable graphs. It is known that for general graphs no efficient distance oracles exist, however, they can be constructed for many classes of graphs. Graphs defined by timetables of regular transport carriers form a specific class which it is not known to admit efficient distance oracles. The thesis should investigate to which extent the known desirable properties (e.g. small highway dimension) are present in these graphs, and/or identify new ones. Analytical study of graph operations and/or experimental verification on real data form two possible approaches to the topic.

Supervisor: doc. RNDr. Rastislav Kráľovič, PhD.
Department: FMFI.KI - Department of Computer Science
Vedúci katedry: doc. RNDr. Daniel Olejár, PhD.
Assigned: 08.11.2011

Approved: 15.11.2011
prof. RNDr. Branislav Rován, PhD.
Guarantor of Study Programme

Student

Supervisor



Univerzita Komenského v Bratislave
Fakulta matematiky, fyziky a informatiky

ZADANIE ZÁVEREČNEJ PRÁCE

Meno a priezvisko študenta: Bc. František Hajnovič
Študijný program: informatika (Jednoodborové štúdium, magisterský II. st., denná forma)
Študijný odbor: 9.2.1. informatika
Typ záverečnej práce: diplomová
Jazyk záverečnej práce: anglický
Sekundárny jazyk: slovenský

Názov: Efektívny výpočet vzdialeností v grafoch spojení lineík.

Cieľ: Cieľom práce je preštudovať možnosti aplikácie výsledkov o distance oracles v grafoch reprezentujúcich dopravné siete na grafy spojení lineík. Otázka, či a aké dôležité vlastnosti ostávajú zachované sa dá riešiť teoreticky pre rôzne triedy grafov a/alebo experimentálne pre reálne dáta.

Vedúci: doc. RNDr. Rastislav Kráľovič, PhD.

Katedra: FMFI.KI - Katedra informatiky

Vedúci katedry: doc. RNDr. Daniel Olejár, PhD.

Dátum zadania: 08.11.2011

Dátum schválenia: 15.11.2011

prof. RNDr. Branislav Rován, PhD.
garant študijného programu

.....
študent

.....
vedúci práce

I hereby declare that I wrote this thesis by myself, only with the help of the referenced literature,
under the careful supervision of my thesis advisor.

.....

Acknowledgements

I would like to thank very much to my supervisor Rastislav Královič for valuable remarks, useful advices and consultations that helped me stay on the right path during my work on this thesis.

I am also grateful for the support of my family during my studies and the work on this thesis.

František Hajnovič

Abstract

In this thesis we deal with queries for optimal connections in timetables on which we have carried out some preprocessing. Based on the analysis of the properties of real-world timetables, we developed exact methods that answer the queries considerably faster than the time-dependent Dijkstra's algorithm implemented with Fibonacci heap priority queue. More specifically, our method *USP-OR-A* with space complexity $\mathcal{O}(n^{1.5})$ achieves average query time $\mathcal{O}(\sqrt{n} \log n)$, outperforming the time-dependent Dijkstra's algorithm up to 7 times in our largest datasets.

Key words: **optimal connection, timetable, Dijkstra's algorithm, Distance oracles, underlying shortest paths**

Abstrakt

V tejto práci sa zaoberáme hľadaním optimálnych spojení v cestovných poriadkoch, na ktorých sme si predpočítali určité informácie. Na základe analýzy reálnych cestovných poriadkov sme vyvinuli exaktné metódy, ktoré na dotaz na optimálne spojenie odpovedajú podstatne rýchlejšie ako časovo závislá implementácia Dijkstrovho algoritmu využívajúca prioritnú frontu na základe Fibonacciho haldy. Presnejšie, náš algoritmus *USP-OR-A* s priestorovou zložitostou $\mathcal{O}(n^{1.5})$ dosahuje časovú zložitosť odpovede na dotaz $\mathcal{O}(\sqrt{n} \log n)$, prekonávajúc časovo závislý Dijkstrov algoritmus takmer 7 krát v našom najväčšom cestovnom poriadku.

Kľúčové slová: **optimálne spojenie, cestovný poriadok, Dijkstrov algoritmus, Dištančné orákulá, podkladové najkratšie cesty**

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Approach	2
1.3	Goals	2
1.4	Organization & conventions	2
2	Preliminaries	3
2.1	Objects	3
2.2	Earliest arrival and optimal connection	6
2.3	(Distance) Oracles	7
2.4	Dijkstra's algorithm	8
3	Related work	9
3.1	Distance oracles and route-planning	9
4	Data & analysis	11
4.1	Data	11
4.2	Analysis of properties	13
5	Underlying shortest paths	16
5.1	<i>USP-OR</i>	17
5.1.1	Analysis of <i>USP-OR</i>	19
5.2	<i>USP-OR-A</i>	21
5.2.1	Analysis of <i>USP-OR-A</i>	23
5.2.2	Correctness of <i>USP-OR-A</i>	26
5.2.3	Modifications of <i>USP-OR-A</i>	27
5.3	Selection of access node set	27
5.3.1	Choosing the optimal access node set	28
5.3.2	Choosing ANs based on node properties	30
5.3.3	Choosing ANs heuristically - the <i>Locsep</i> algorithm	31
5.4	Performance and comparisons	34
5.4.1	Performance of <i>USP-OR</i>	34
5.4.2	<i>USP-OR-A</i> with <i>Locsep</i>	37
5.4.3	<i>USP-OR-A</i> with <i>Locsep Max</i>	40
6	Neural network approach	42
7	Application TTBlazer	43
8	Conclusion	44
	Appendix A File formats	45

1 Introduction

World is getting smaller every day, as new technologies constantly make communication and travelling faster and more effective than yesterday. Road network, Internet and many other networks are becoming more evolved and denser which also brings along new problems. In order to fully take advantage of such huge networks, we must have efficient algorithms that operate on these networks and give us answers to many questions. Among many others, one that we take particular interest in is the question: “What is the shortest path from place x to place y ”?

In different networks, this question can make different sense. In the road network, we would like to obtain a sequence of intersections we have to go through in order to reach our destination, driving the shortest possible time (or the smallest possible distance). GPS devices and the likes of Google maps have to deal with this problem. In case of the Internet network, we might be interested in the shortest path to a destination computer in terms of router hops. In a network of social acquaintances, the smallest number of persons connecting us e.g. with guitarist Mark Knopfler or Liona Boyd could be expressed as a shortest path problem. Many problems in artificial intelligence (e.g. planning of actions) can be expressed, or include, looking for shortest paths.

The tremendous amount of work done in this area signifies the importance of quick distance or shortest path retrieval in graphs. A simple Dijkstra’s or A* algorithm no longer comply to the requirements of today’s applications, in which a server often has to answer hundreds of shortest path queries per second in a large-scale networks. To speed up the mentioned algorithms we usually sacrifice generality and concentrate on a particular type of network, if not only on one concrete network.

In this thesis, the type of network we deal with is the one representing timetable connections, where nodes are the stations and arcs represent a direct connection between the two stations. We will talk in more details about this in following sections. However, this network has one substantial difference that we would like to point out - it is time-dependent. That means that the shortest path from station x to station y may have different solutions depending on the time when we start at station x . Therefore, we will not talk about shortest paths and distances, but rather about optimal connections and earliest arrivals and each query will now bear a third parameter - the departure time from x .

To informally develop the discussion about optimal connections in timetables, we will now clarify the motivation, approach and the goals of this thesis.

1.1 Motivation

We have already sketched-out the motivation in the introductory text. We consider that a server (hosting e.g. journey-planning application) has to answer many queries per given time unit. What does it mean many? British National Rails Enquiries website that hosts journey planner supports over 1 million queries per day [?]. Even if these queries were distributed evenly throughout the whole day, there would still be more than 11 queries per second. That is why the search engine run for each query has to be fast enough to provide an answer.

11 queries per second is probably not a big issue. There is about 2500 railway stations in Great Britain and a current state of the art computer with basic implementation of a time-dependent Dijkstra’s algorithm (to be talked about later) would be able handle the mentioned load without any problems. However, things get more difficult on a bigger scale, in rush hours and when additional

requirements are posed on the search results (transfers, cost of travel or simply outputting more results the user can choose from).

In shortest path routing on road networks very much has been done to speed-up the query times using pre-processing on the input graph (for a good review such methods, see [?]). Some developed methods answer distance queries ≈ 1000000 faster than the Dijkstra's algorithm. The timetable scenario has so far seen much smaller speed-ups, one reason for this being that the adaptation of the many techniques used for road networks to the time-dependent scenario was not so straightforward [?].

1.2 Approach

We have mentioned that to get more effective algorithms with better query times, we need to focus on a special type of network and take advantage of its properties. Another thing we can do in this case is to pre-compute some information on the particular timetable and to use this information later to speed-up the answering of the queries. This is not a new technique and in the shortest path routing it is commonly referred to as creating distance oracle [?]. Our approach is analogical - we simply deal with graphs representing the timetables ¹ and look for optimal connections. We will go more into the details in the preliminaries section 2.

1.3 Goals

We have set two main goals for this thesis:

- **Analyse real-world timetables** and their properties. More specifically, we were interested in the sparsity of the underlying graphs, their connectivity, average and maximal degrees, average optimal connection sizes, betweenness centrality distribution, highway dimension... We will talk about these mostly in the section 4.
- **Develop methods with fast query times** for optimal connections, based on the approach we have sketched in the previous subsection. For this purpose, we use also the outcomes of our analysis.

1.4 Organization & conventions

¹Hence the name of this thesis - Distance oracles for timetable graphs, the “distance” being part of the title because the term “distance oracle” is generally recognized.

2 Preliminaries

In this section, we provide most of the definitions and terminology used throughout the thesis.

2.1 Objects

First, we will formalize the notion of a timetable and its derived graph forms, the underlying graph and terms related to these objects.

Definition 2.1. Timetable (TT)

A timetable is a set $T = \{(x, y, p, q) \mid p, q \in \mathbb{N}, p < q\}$.

- Elements of T (the 4-tuples) are called **elementary connections**. For an elementary connection $e = (x, y, p, q)$:
 - $from(e) = x$ is the **departure city**
 - $to(e) = y$ is the **arrival/destination city**
 - $dep(e) = p$ is the **departure time**
 - $arr(e) = q$ is the **arrival time**
- The set of all **cities** will be denoted as $ct_T = \{x \mid (x, y, p, q) \in T \text{ or } (y, x, p, q) \in T\}$ and the number of cities as n_T
- Pairs (x, p) or (y, q) such that $(x, y, p, q) \in T$ form the set of **events** ev_T . The set of events in a specific city x is $ev_T(x) = \{(x, t) \mid (x, y, t, q) \in T \text{ or } (y, x, p, t) \in T\}$
- Let $tlow_T = \min_{e \in T} dep(e)$ and $thigh_T = \max_{e \in T} arr(e)$. The value $r_T = thigh_T - tlow_T$ is called the **time range** of the timetable.
- **Height** of the timetable is the maximum number of events in a city: $h_T = \max_{x \in cities_T} \{|ev_T(x)|\}$

Let us describe some the defined terms more informally. An elementary connection corresponds to moving from one stop to the next one, e.g. with a bus (thus we disregard the notion of *lines*, i.e. getting on and off). Note that we express time as an integer - throughout this paper, this integer will represent the minutes elapsed from the time 00:00 of the first day. Thus we may take the liberty of talking about time in integer or *days hh:mm* format, as convenient at the moment. Lastly, an event simply represent an arrival or departure of a e.g. train at some station. The remaining terms should be clear enough.

Place		Time	
From	To	Departure	Arrival
A	B	10:00	10:45
B	C	11:00	11:30
B	C	11:30	12:10
B	A	11:20	12:30
C	A	11:45	12:15

Table 2.1: An example of a timetable - the set of elementary connections (between pairs of **cities**). An example of an event is a pair (A, 10:00), when some el. connection departs from A.

Following is a definition of a connection.

Definition 2.2. Connection

A connection from a to b is a sequence of elementary connections $c = (e_1, e_2, \dots, e_k), k \geq 1$, such that $from(e_1) = a$, $to(e_k) = b$ and $\forall i \in \{2, \dots, k\} : (to(e_i) = from(e_{i-1}), arr(e_i) \geq dep(e_{i-1}))$.

- Connection **starts** at the departure time $\text{start}(\mathbf{c}) = \text{dep}(e_1)$ and **ends** at the arrival time $\text{end}(\mathbf{c}) = \text{arr}(e_k)$.
- We also extend $\text{from}(\mathbf{c}) = \text{from}(e_1)$ and $\text{to}(\mathbf{c}) = \text{to}(e_k)$
- **Length** of the connection is $\text{len}(\mathbf{c}) = \text{end}(\mathbf{c}) - \text{start}(\mathbf{c})$
- **Size** of the connection is $\text{size}(\mathbf{c}) = k^2$
- We will denote the set of **all connections** from a to b in a timetable T as $\mathbf{C}_T(a, b)$. We also define $\mathbf{C}_T = \cup_{a,b} \mathbf{C}_T(a, b)$

So we understand connection as a (valid) sequence of elementary connections.

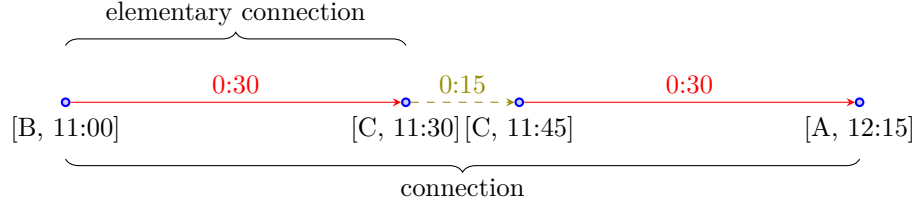


Figure 2.1: A valid connection made out of **elementary connections** (and **waiting**, which is implicit).

Next, we continue with the underlying graph - a graph representing basically the map on top of which the timetable operates.

Definition 2.3. Underlying graph (UG graph)

The underlying graph of a timetable T , denoted \mathbf{ug}_T , is an oriented graph (V, E) , where V is the set of all timetable cities and $E = \{(x, y) \mid \exists (x, y, p, q) \in T\}$

- By \mathbf{m}_T we will denote the number of arcs in the UG

Note, that we do not specify the weights of the edges in the underlying graph - they will be specified based on the current usage of the UG. Most of the time, however, if we work with a weighted UG, the weight of an arc will be the length of the shortest elementary connection on that arc. More specifically, $w(x, y) = \min_{(x, y, p, q) \in T} (q - p) \ \forall (x, y) \in E(\mathbf{ug}_T)$. Such weighted UG will be called **optimistic** (denoted $\mathbf{ug}_T^{\text{opt}}$).

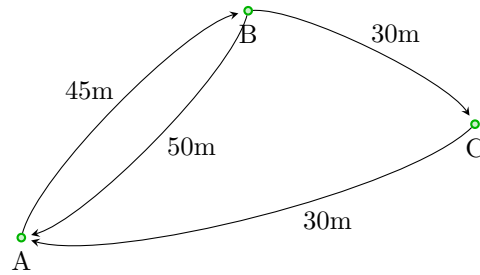


Figure 2.2: An optimistic underlying graph of the timetable 2.1. The nodes are the **cities** of the timetable.

²We will use similar terminology when talking about paths - the *size* is the number of vertices (hops) in the path while the *length* refers to the actual distance (sum of weights of the edges in the path).

If we want to represent the timetable by a graph, there are two most common options [?] - the time-expanded and time-dependent graph.

Definition 2.4. Time-expanded graph (TE graph)

Let T be a timetable. Time-expanded graph from T , denoted te_T , is an oriented graph (V, E) whose vertices correspond to events of T , that is $V = \{(x, t) \mid (x, t) \in ev_T\}$. The edges of G are of two types

1. $([x, p], [y, q]) \forall (x, y, p, q) \in T$ - the so called **connection edges**
2. $([x, p], [x, q]) [x, p], [x, q] \in V, p < q$ and $\nexists [x, r] \in V: p < r < q$. - the so called **waiting edges**

Weight of the edge $([x, p], [y, q])$ is $w([x, p], [y, q]) = q - p$.

Informally, an edge in TE graph represent either the travelling with an elementary connection or waiting for the next event in the same city. Also, the time range and height of a timetable could be easily illustrated on the TE graph (see picture 2.3).

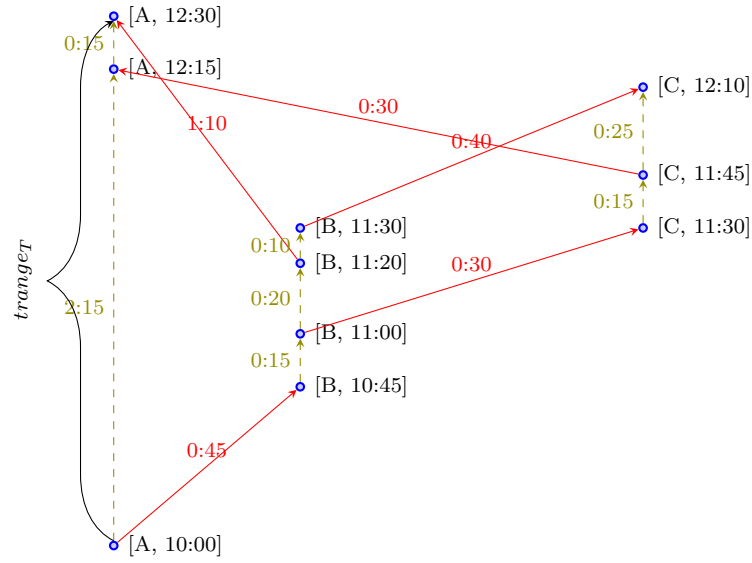


Figure 2.3: Time-expanded graph of the timetable 2.1. Nodes represent the **events**. There are **connection** and **waiting** edges (dashed). The time range is 2h:30m and the height is 4 (as there are 4 events in city B).

Definition 2.5. Time-dependent graph (TD graph)

Let T be a timetable. Time-dependent graph from T , denoted td_T , is an oriented graph (V, E) whose vertices are the timetable cities and $E = \{(x, y) \mid \exists (x, y, p, q) \in T\}$. Furthermore, the weight of an edge $(x, y) \in E$ is a piece-wise linear function $w(x, y) = f_{x,y}(t) = q - t$ where q is:

- $\min\{arr(e) \mid e \in T, dep(e) \geq t\}$
- ∞ , if $dep(e) < t \forall e \in T$

Intuitively, the TD graph is simply the UG graph where each arc carries a function specifying the traversal time of that arc at any time. For an example, see picture 2.5: The latest point of every linear segment is called the **interpolation point** and it corresponds to an elementary connection (its coordinates are $dep(e), len(e)$ for corresponding el. connection e). Note that a list of all interpolation points fully defines the piece-wise linear function.

The algorithms in this thesis use almost exclusively the TD graphs, mainly because they are less space consuming. Also, time-dependent Dijkstra searches are a bit faster on TD graphs,

because the search space that has to be explored is smaller. On the other hand, TE graphs are more flexible when we need to take additional search parameters into consideration (like transfers, travel costs). Since we will not talk about these, TD graphs are more suitable.

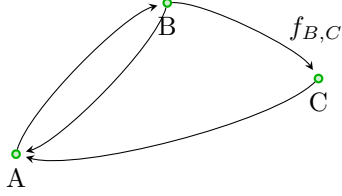


Figure 2.4: Time-dependent graph of the timetable 2.1. The nodes are the cities.

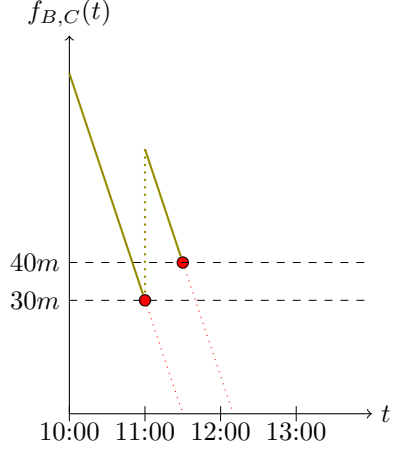


Figure 2.5: Piece-wise linear function - traversal times for the arc (B, C) . The highlighted points are the interpolation points.

To sum up, there are four main types of objects we will be working with:

- Timetable (TT)
- Underlying graph (UG)
- Time-expanded graph (TE)
- Time-dependent graph (TD)

For further reference, we will call **timetable objects** those, that fully represent a timetable (TT, TE, TD) and **graph objects** those, that can be viewed as a graph (UG, TE, TD).

Note: Throughout this paper, we will relax a bit the notation and leave out subscripts (e.g. $ug_T \rightarrow ug$, $n_T \rightarrow n$, etc.) in situations, where the context is clear enough.

2.2 Earliest arrival and optimal connection

Now we would like to formulate the main problems this thesis deals with.

Definition 2.6. Earliest arrival problem (EAP)

Given a timetable T , departure city x , destination city y and a departure time t , the task is to determine $\mathbf{t}_{(x,t,y)}^* = \min_{c \in C_T(x,y)} \{t + \text{len}(c) \mid \text{start}(c) \geq t\}$.

- We will refer to the tuple (x, t, y) as an **EAP instance**, or an **EAP query**
- The time $\mathbf{t}_{(x,t,y)}^*$ is called the **earliest arrival (EA)** for the given EAP instance

A bit more difficult version of this problem is one, where we require to actually output the connection ending at time given by EA.

Definition 2.7. Optimal connection problem (OCP)

Given a timetable T , departure city x , destination city y and a departure time t , the task is to determine the **optimal connection (OC)** $\mathbf{c}_{(a,t,b)}^* = \text{argmin}_{c \in C_T(a,b)} \{t + \text{len}(c) \mid \text{start}(c) \geq t\}$.

The instance/query in case of the optimal connection problem has the same form as EAP query. Also, note that the OCP is at least as hard to solve as EAP since having the optimal connection implies the optimal (earliest) arrival time.. In order to avoid technical issues in later parts of the thesis, we will assume the optimal connection is unique (i.e., there is not a different connection with the same end time) or that ties are won by a lexicographically first connection.

Example 2.1. Consider our timetable from table 2.1. For the EAP instance $(B, 10:45, A)$, the earliest arrival (EA) is 12:15 and the optimal connection (OC) is $((B, C, 11:00, 11:30), (C, A, 11:45, 12:15))$, as could be easily seen from picture 2.6 of the TE graph.

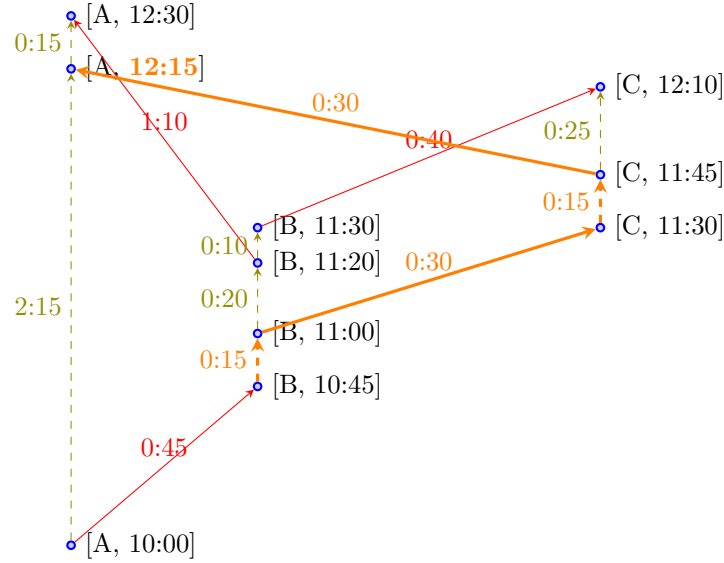


Figure 2.6: Optimal connection and earliest arrival time are marked in **bold**.

2.3 (Distance) Oracles

The term *distance oracle* was first coined in 2001 by Thorup and Zwick [?], when talking about quick shortest path (or distance) computations on graphs. One approach to this problem is to pre-compute some information on the graph to speed-up answering of the queries. The paper of Thorup and Zwick was dealing with trade-offs among the time complexity of the pre-computation, the amount of pre-computed information, the speed-up in query times and the accuracy of the answers. Since the pre-computed data structure is something that helps us answer the queries more efficiently, it resembles an oracle, thus the term distance oracle.

In this thesis, we will discuss methods that behave the same way, but deal with the earliest arrival problem (or optimal connection problem) - there is some pre-processing of the timetable with a resulting data structure that speeds up answering subsequent queries. To formalize this a little more, we will refer to this kind of methods as **oracle based methods**. For such a method m , we are interested mainly in its four parameters:

- **Preprocessing time** ($prep(m)$) - the time complexity of the pre-computation
- **Preprocessed space** ($size(m)$) - the space complexity of the pre-computed data structure (the so called **oracle**)

- **Query time** ($qtime(m)$) - the time complexity of answering a single query
- **Stretch** ($stretch(m)$) - the worst-case ratio against the optimal value of earliest arrival (the lower, the better)

The preprocessing time is probably the least critical resource. A reasonable polynomial should bind its time complexity, depending on the computational power of the user and the scale of the timetable. The size of the preprocessed oracle is much more important - in the optimal case, it should be bound by the space complexity of the timetable itself. Optimality of the query time depends on which problem we are solving. If we query for the whole optimal connection, we have to count with a time complexity at least proportional to the diameter of the underlying graph (as connections could be that long, or even longer). If we require only the EA value as an output, much better speed-ups could be expected. The stretch should be of course as low as possible.

2.4 Dijkstra's algorithm

Throughout this thesis, we will often use Dijkstra's algorithm and its modifications both as a part of our algorithms and as a reference point against which we will compare the performance of our methods. This is a common practice. Researchers working on methods answering distance or shortest path queries in road networks commonly use the term *speed-up*, i.e. *how many times faster* is their algorithm against the Dijkstra's algorithm.

Dijkstra's algorithm is originally an algorithm that looks for shortest paths in weighted oriented graphs. It was published by E. W. Dijkstra in 1959 [?] and we will not explain it at this place, as the algorithm is very well explained at many other places (e.g. [?]). For a good summary of Dijkstra's algorithm related implementations and publications see [?].

As our task is to compute earliest arrivals or optimal connections instead of distances and shortest paths, our "reference point" will be a slightly modified Dijkstra's algorithm called **time-dependent Dijkstra's algorithm** ?? (or TD Dijkstra for short). The algorithm is run on a time-dependent graph and works just like the ordinary Dijkstra's algorithm, except that the weight of each arc (x, y) is determined for the time t at which we had settled vertex x .

If we assume that the evaluation of an arc by the cost function of the TD graph is implemented in constant time, the running time of the TD Dijkstra is $\mathcal{O}(n^2)$, just like the normal Dijkstra's algorithm. On sparse graphs, this bound can be improved using a quick data structure to determine the next node we settle. A good option is a priority queue implemented as a *Fibonacci heap*, which implements deletion in $\mathcal{O}(\log n)$ and all other operations in constant amortized time [?]. This yields the running time of TD Dijkstra $\mathcal{O}(n \log n + m)$.

We may therefore introduce a fifth parameter of our oracle based methods, the speed-up:

Definition 2.8. *Speed-up* ($spd(m)$)

A *speed-up* of an oracle based method m is the ratio $\frac{qtime_{avg}(TDDijkstra)}{qtime_{avg}(m)}$ where $qtime_{avg}(m')$ is the average query time of the respective oracle based method m' ³.

The definition is rather loose in the sense that we may refer to a concrete speed-up of the method on a concrete dataset or a general, theoretical speed-up expressed as a function of the size of input.

³Note that we may also consider the TD Dijkstra algorithm to be an oracle based method - it just happens that it does not require any preprocessing.

3 Related work

In this section, we summarize the work related to the subject of this thesis. Apart from the papers discussing searching for optimal connections and earliest arrivals in time-dependent scenarios, we also briefly summarize the research done on route planning in road networks and on distance oracles in general.

3.1 Distance oracles and route-planning

We have already mentioned in section 2 the paper of Thorup and Zwick [?] where the term “distance oracle” originated. The authors have shown, that given an undirected weighted graph of n vertices and m edges and a chosen integer $k \geq 1$, we can build a distance oracle such that:

- preprocessing takes $O(kmn^{1/k})$ expected time
- resulting distance oracle is of size $O(kn^{1+1/k})$
- answering queries takes $O(k)$ time
- stretch is at most $2k - 1$

Moreover, the authors have reasoned that their construction is essentially optimal with respect to space - i.e., if we want to have exact and constant-time answers, we will in general be forced to pre-compute $\Omega(n^2)$ information. The parameter k however provides a nice option to make trade-offs between the four parameters, as depicted on figure 3.1.

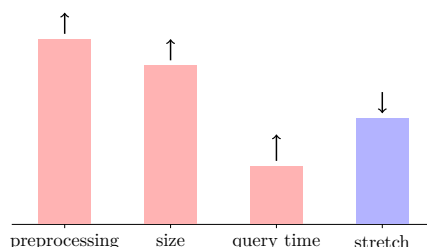


Figure 3.1: By moving k , we can achieve compromises between the four parameters of the distance oracle.

Another work by Gavaille et al. [?] concerned distance labelling - a somewhat restricted version of a distance oracle where we assign each node in the graph its distance label. This is again only some pre-computed information and upon a query from x to y , we should be able to figure out their distance only using the corresponding distance labels. In the paper it is shown that for all n , there exist infinitely many graphs of n vertices for which we have an exact distance labelling scheme of a small overall size ($\mathcal{O}(n \log n)$), but for which the process of figuring out the distance from the labels takes too long from practical point of view.

Even though these results imply that we cannot create a sufficiently small efficient distance oracle in general, it may still be possible for sub-classes of general graphs, or even better, for a single particular graph. In that respect, the road network is the point of interest and fortunately it has a few “nice” properties (it is sparse, almost planar, the maximum degree of the node is small...) which

made it possible to design exact and efficient algorithms with extremely fast query times. To name a few of these:

- Highway hierarchies (2005, [?])
- Transit node routing (2006, [?])
- Contraction hierarchies (2008, [?])

A very good summary of the techniques devised for road network route planning up to 2009 can be found in [?].

The work [?] gives an exhaustive and comprehensive discussion regarding shortest path queries in general, while suggesting efficient distance oracle for power-law graphs.

4 Data & analysis

In this section we would like to introduce the timetable datasets we were working with and provide the results of the analysis which we carried out on the data. The main reason for this analysis is that it gives some insight into the properties of the timetables, and thus may contribute to the make an oracle based method with better qualities.

4.1 Data

We have obtained timetable datasets from numerous sources, in varying formats and of different types. Some of them were freely available on the Internet while others were provided by companies upon demand. Let us briefly describe each of these timetables.

The dataset *air01* contains schedules of **domestic flights in United States** for the January of 2008. It is not comprehensive in the sense that it contains entries only for flights of some of the major airports in US. However it is large enough for our purposes (almost 300 airports). This dataset is just a fraction of the data that are freely available at the pages of American Statistical Association ⁴ in CSV format.

Timetable *cp.sk* represent the **regional bus** schedules from the areas of **Ružomberok and Žilina, Slovakia**. The data were provided by the company in charge of the *cp.sk* portal - Inprop s.r.o. . The timetable concern about 1900 bus stops and came in a JDF 1.9 format ⁵. Apart from the actual schedules, the data in JDF contain numerous other information, which were not relevant for our purposes. From both timetables, we have extracted subsets with a time range of one day.

The *gb-coach* and *gb-train* timetables are freely available from National Public Transport Data Repository (NPTDR) at <http://data.gov.uk> in an ATCO-CIF format. They are not actually timetables but rather weekly snapshots of national public transport journeys made by coach and train in Great Britain. The datasets contain about 2500 stations each.

The *montr* dataset is part of a public feed for **Greater Montreal public transportation**, available at Google Transit Feeds ⁶. The data are in a GTFS format (defines relations between CSV files listing stations, routes, stop-times...) and were made available by Montreal's Agence métropolitaine de transport. Our timetable *montr* corresponds to daily schedules of the Chambly-Richelieu-Carignan bus services (more than 200 bus stops).

Also in GTFS format come the data of **French railways** operated by company SNCF, publicly available at their website ⁷. The schedules are weekly and there were two of them: one for intercity trains and one for TER trains (regional trains). Thus the three timetables *snCF-inter* (366 stations), *snCF-ter* (2637 stations) and their union *snCF* (2646 stations).

Finally, one more country-wide railway timetable was provided by ŽSR, the company in charge of the **Slovak national railways**. This timetable was exported in a MERITS format and its time range is for one year. The number of stations in *zsr* dataset is 233.

With the help of Python and Bash scripts, we converted each of these datasets to our timetable format (described in appendix A). This timetables were then loaded by our application TTBlazer, which can further generate sub-timetables (with less stations or smaller time range), underlying graphs and TE and TD graphs.

⁴<http://stat-computing.org/dataexpo/2009/the-data.html>

⁵Jednotný datový formát (JDF).

⁶<http://code.google.com/p/googletransitdatafeed/wiki/PublicFeeds>

⁷<http://test.data-sncf.com/index.php/ter.html>

For a summary of the used timetables’ descriptions, see table 4.1 and for their main properties, refer to table 4.2.

Name	Description	Format	Provided by	Publicly available
<i>air01</i>	domestic flights (US)	CSV	American Stat. Assoc.	✓
<i>cpsk</i>	regional bus (Ružomberok & Žilina, SVK)	JDF 1.9	Inprop s.r.o.	✗
<i>gb-coach</i>	country-wide buses (GB)	ATCO-CIF	NPTDR	✓
<i>gb-train</i>	country-wide rails (GB)	ATCO-CIF	NPTDR	✓
<i>montr</i>	public transport (Montreal, CA)	GTFS	Montreal AMT	✓
<i>sncf</i>	country-wide rails (FRA)	GTFS	SNCF	✓
<i>zsr</i>	country-wide rails (SVK)	MERITS	ŽSR	✗

Table 4.1: Datasets descriptions.

Name	El. conns.	Cities	UG arcs	Time range	Height
<i>air01</i>	601489	287	4668	1 month	24374
<i>cpsk</i>	97916	1905	5093	1 day	370
<i>gb-coach</i>	260710	2448	5793	1 week	3140
<i>gb-train</i>	1714535	2555	8335	1 week	7978
<i>montr</i>	7153	217	349	1 day	363
<i>sncf</i>	416302	2646	7994	1 week	2679
<i>sncf-inter</i>	22750	366	901	1 week	1052
<i>sncf-ter</i>	393587	2637	7647	1 week	2646
<i>zsr</i>	932052	233	588	1 year	60308

Table 4.2: Main properties of the timetables. The value of time range is approximate.

To see better the differences in the properties of different timetable types (train, flight, bus...), we made sub-timetables with 200 cities and with the upper bound on time range being 1 day ⁸ ($thigh_T < 1 \text{ day } \forall T$) from each of our dataset. We name these datasets by appending to the original name “-200d” ⁹. See table 4.3 for details.

Name	El. conns.	Cities	UG arcs	Height
<i>air01-200d</i>	19010	200	3973	772
<i>cpsk-200d</i>	14747	200	592	370
<i>gb-coach-200d</i>	2760	200	564	498
<i>gb-train-200d</i>	24323	200	792	957
<i>montr-200d</i>	6841	200	320	355
<i>sncf-200d</i>	4192	200	611	269
<i>sncf-inter-200d</i>	2172	200	493	128
<i>sncf-ter-200d</i>	8469	200	600	419
<i>zsr-200d</i>	2031	200	454	133

Table 4.3: 200-station sub-timetables with the time range of one day.

Also, to further justify our choice of using TD graphs instead of TE graphs in this thesis, we provide their space consumption comparison in table 4.4.

⁸We took all elementary connections that were within our time range. From this timetable, we made an UG and its (random) sub-graph of 200 cities. Finally we selected only those elementary connections, that were on top of this sub-graph to form a timetable with 200 cities and the desired (maximal) time range.

⁹Similarly, “-d” would mean “with daily time range”.

Name	TD graph			TE graph		
	Nodes	Arcs	Size (MB)	Nodes	Arcs	Size (MB)
<i>air01</i>	287	4668	27	715211	1307432	72
<i>cpsk</i>	1905	5093	5	95601	189205	11
<i>gb-coach</i>	2448	5793	12	259589	512862	32
<i>gb-train</i>	2555	8335	79	2042316	3745751	263
<i>montr</i>	217	349	0.4	7182	13992	0.9
<i>sncf</i>	2646	7994	19	758867	1166646	85
<i>sncf-inter</i>	366	901	1.1	39765	60602	4.6
<i>sncf-ter</i>	2637	7647	18	720651	1107301	81
<i>zsr</i>	233	588	42	1706077	2637896	173

Table 4.4: Space consumption of time-dependent vs. time-expanded model. The number of nodes and arcs for TD graph is the same as for the corresponding underlying graph.

4.2 Analysis of properties

First we will take a look at the optimal connection *sizes* (size is the number of el. connections) in the timetables. For a given timetable T , we will denote the average optimal connection size as γ_T and will call it the **optimal connection radius** (OC radius). We computed an approximate OC radius for each of our datasets by measuring an average connection size of sufficiently many OCs. The results in table 4.5 indicate that the average OC size generally falls under \sqrt{n} .

Next we would like to get an idea of the sparsity of the underlying graphs. We see from the table 4.2 that the graphs are pretty sparse (again, with exception of *air01*), but we would like to make sure that the sparsity is uniform. More specifically, we will be interested in the δ -density:

Definition 4.1. δ -density

A graph G of n vertices and m arcs is δ -dense $\iff \forall G' \subseteq G, n' \geq \sqrt[4]{n} : \frac{m'}{n'} \leq \delta$

- For a timetable T , we will denote its **density** parameter¹⁰ as $\delta_T = \min\{\delta \mid ug_T \text{ is } \delta\text{-dense}\}$

To find out at least approximate δ_T values for our timetables, we have randomly sampled their UGs for (connected) sub-graphs of various sizes (starting from $\sqrt[4]{n}$ ¹¹). In table 4.6 you can see the maximal density found during the sampling.

The density is related to the **average degree** deg_{avg} in the graph, since in oriented graphs:

$$deg_{avg} = \frac{m}{n}$$

So the average degree is a lower bound on the graph's density. Table 4.7 lists the average and maximal degrees in the underlying graphs.

We would also assume, that the underlying graphs of each timetable will be **connected** (and even strongly connected), or at least that the largest connected component spans almost the whole graph. From the table 4.8 we may see that this assumption holds.

¹⁰Note that this has nothing to do with the frequency of elementary connections.

¹¹The choice of $\sqrt[4]{n}$ will be justified later, during the analysis of the algorithms.

Name	γ_T	Max. OC size found	\sqrt{n}
<i>air01</i>	2.4	8	16.9
<i>cpsk</i>	40.8	162	43.6
<i>gb-coach</i>	25.2	128	49.5
<i>gb-train</i>	25.6	111	50.5
<i>montr</i>	21.1	63	14.7
<i>sncf</i>	36.8	111	51.4
<i>sncf-inter</i>	17.1	58	19.1
<i>sncf-ter</i>	48.0	167	51.3
<i>zsr</i>	15.0	57	15.3

Table 4.5: With one exception, OC radius is less than \sqrt{n} (this was expected, as *montr* is the only timetable with “geographically one dimension long” - all other timetables span areas with more uniform shape). Note extremely low value for airline timetable - this is due to the fact that UGs of airline timetables have small-world characteristics [?]. Another thing we may notice is that regional timetables (*cpsk*, *sncf-ter*) have higher OC radius than country-wide and inter-city timetables. We also point out that the inter-city trains in French railways decrease the average optimal connection size by one about third.

Name	Maximal δ_T found
<i>air01</i>	34.5
<i>cpsk</i>	4.1
<i>gb-coach</i>	5.0
<i>gb-train</i>	5.8
<i>montr</i>	1.9
<i>sncf</i>	5.0
<i>sncf-inter</i>	3.0
<i>sncf-ter</i>	4.8
<i>zsr</i>	3.2

Table 4.6: Approximate density of the underlying graphs.

Name	Avg. degree	Max. degree
<i>air01</i>	16.3	166
<i>cpsk</i>	2.7	27
<i>gb-coach</i>	2.4	103
<i>gb-train</i>	3.3	30
<i>montr</i>	1.6	5
<i>sncf</i>	3.0	27
<i>sncf-inter</i>	2.5	12
<i>sncf-ter</i>	2.9	27
<i>zsr</i>	2.5	12

Table 4.7: Average and maximal degree in the underlying graphs.

Name	n	Connectivity		Strong connectivity	
		Connected	Largest comp.	Connected	Largest comp.
<i>air01</i>	287	✓	287	✗	286
<i>cpsk</i>	1905	✓	1905	✗	1903
<i>gb-coach</i>	2448	✗	2374	✗	2332
<i>gb-train</i>	2555	✓	2555	✓	2555
<i>montr</i>	217	✗	211	✗	209
<i>sncf</i>	2646	✓	2646	✗	2594
<i>sncf-inter</i>	366	✗	328	✗	316
<i>sncf-ter</i>	2637	✓	2637	✗	2583
<i>zsr</i>	233	✓	233	✗	225

Table 4.8: Connectivity of underlying graphs.

In the previous section 3 we have mentioned the highway dimension [?] as a parameter which, when being low, guarantees low query times for certain route-planning methods. Here we were interested in the highway dimension of our underlying graphs.

Definition 4.2. Highway dimension

Highway dimension $HD(G)$ for a directed, edge-weighted graph $G = (V, E)$ is the smallest integer h , such that:

$$\forall r \in R^+, \forall u \in V, \exists S \subseteq B_{u,2r}, |S| \leq h, \forall v, w \in B_{u,2r}: \\ \text{if } r < |P(v, w)| \leq 2r \text{ and } P(v, w) \subseteq B_{u,2r} \text{ then } P(v, w) \cap S \neq \emptyset$$

where:

- $P(v, w)$ is the **shortest path** between v and w
- $B_{u,r} = \{v \in V \mid |P(u, v)| \leq r \text{ or } |P(v, u)| \leq r\}$ and is called **ball** of radius r centred at u .

Intuitively, a graph has a low HD, if for any r we have a *sparse* set of vertices S_r , such that every shortest path longer than r includes a vertex from S_r . By the set being sparse, we mean that every ball of radius $\mathcal{O}(r)$ contains just a few elements of S_r .

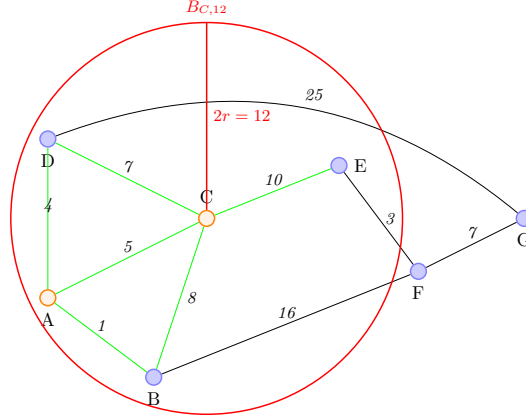


Figure 4.1: Demonstration of a definition of HD. We chose some r ($r = 6$) and some vertex v ($v = C$) to root the ball $B_{v,2r}$. All the shortest paths *longer* than r *inside* the ball have to contain a vertex from S (orange vertices C and A in our case). The upper bound on $|S|$, considering any ball with any radius, is the required highway dimension. Note: in our case, we had to choose also A to set S , since a shortest path from B to D does not include C .

5 Underlying shortest paths

In section 2 we have defined a timetable as a set of elementary connections. While do not pose any other restrictions on this set or on the elementary connections themselves, the real world timetables usually have a specific nature. Quite often are the connections repetitive, that is, the same sequence of elementary connections is repeated in several different moments throughout the day.

Another thing we may notice is that if we talk about *optimal* connections between a pair of distant cities u and v , we are often left with a few possibilities as to *which way should we go*. This is not only because the underlying graph is usually quite sparse¹², but also because for longer distances we generally need to make use of some express connection that stops only in (small number of) bigger cities.

Thus the main idea which will repeat often throughout this section: *when carrying out an optimal connection between a pair of cities, one often goes along the same path regardless of the starting time*.

To formalize this idea, we will introduce the definition of an *underlying shortest path* - a path in UG that corresponds to some optimal connection in the timetable. To do this, we will first define a function *path* that extracts the **underlying path** (trajectory in the UG) from a given connection. Let c be a connection $c = (e_1, e_2, \dots, e_k)$.

$$\mathbf{path}(c) = \mathit{shrink}(\mathit{from}(e_1), \mathit{from}(e_2), \dots, \mathit{from}(e_k), \mathit{to}(e_k))$$

Note, that if the connection involves waiting in a city (as e.g. in picture 5.1), $e_x^i = e_x^{i+1}$ for some i . That is why we apply the *shrink* function, which replaces any sub-sequences of the type (z, z, \dots, z) by (z) in a sequence. This was rather technical way of expressing a simple intuition - for a given connection, the *path* function simply outputs a sequence of visited cities. Now we can formalize the underlying shortest path.

Definition 5.1. Underlying shortest path (USP)

A path $p = (v_1, v_2, \dots, v_k)$ in UG_T is an **underlying shortest path** if and only if $\exists t \in \mathbb{N} : p = \mathbf{path}(c_{(v_1, t, v_k)}^*), c_{(v_1, t, v_k)}^* \in C_T$

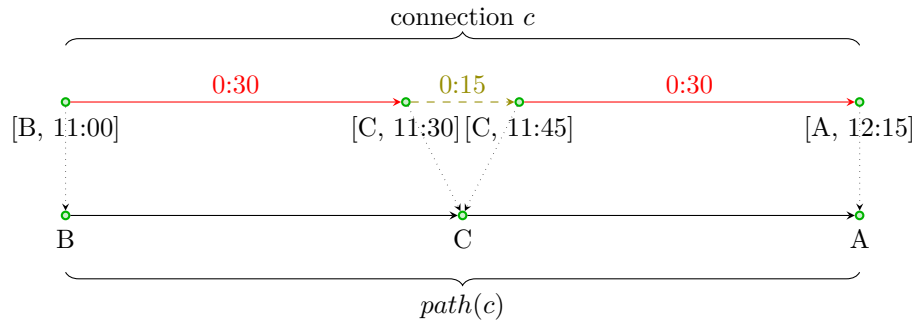


Figure 5.1: The *path* function applied on a connection to get the underlying path.

Please note that the terminology might be a bit misleading - an USP is not necessarily a shortest path in the given UG. Connections on a shortest path may simple require too much waiting (the

¹²Maybe with exception of the airline timetables, which tend to be more dense.

el. connections simply do not follow well enough one another) and thus it might be that travelling along the paths with greater distance proof to be faster options.

5.1 USP-OR

We can easily extract the underlying path from a given connection. Now let us look at this from the other way - if, for a given EA query, we know the underlying shortest path, can we reconstruct the optimal connection? One thing we could do is to blindly follow the USP and at each stop take the first elementary connection to the next stop on the USP. This simple algorithm called *Expand* is described in algorithm 5.1.

Algorithm 5.1 Expand

Input

- timetable T
- path $p = (v_1, v_2, \dots, v_k)$, $v_i \in ct_T$
- departure time t

Algorithm

```

 $c$  = empty connection
 $t' = t$ 
for all  $i \in \{1, \dots, k - 1\}$  do
     $e = \operatorname{argmin}_{e' \in C_T(v_i, v_{i+1})} \{dep(e') \mid dep(e') \geq t'\}$       # take first available el. conn.
     $t' = \operatorname{arr}(e)$ 
     $c := e$       # add the el.conn to the resulting connection
end for

```

Output

- connection c
-

Will we get an optimal connection if we expanded all possible USPs between a pair of cities? We show that we will, provided the timetable has no *overtaking* [?] [?] of elementary connections.

Definition 5.2. Overtaking

An elementary connection e_1 **overtakes** e_2 if, and only if $dep(e_1) > dep(e_2)$ and $arr(e_1) < arr(e_2)$.

Lemma 5.1. Let T be a timetable without overtaking, (x, t, y) an EA query in this timetable and $\mathcal{P} = \{p_1, p_2, \dots, p_k\}$ a set of all USPs from x to y . Define $c_i = \operatorname{Expand}(T, p_i, t)$ to be the connection returned by the algorithm Expand 5.1. Then $\exists j : c_j = c_{x,t,y}^*$.

Proof. The optimal connection $c_{x,t,y}^*$ has an USP p which must be present in the set \mathcal{P} , as it is the set of all USPs from x to y . So $p = p_j = (v_1, v_2, \dots, v_l)$ from some j . We want to show that c_j is the optimal connection. This may be shown inductively:

1. *Base:* Expand reaches city $v_1 = x$ as soon as possible (since the connection just starts there)
2. *Induction:* Expand reached city v_i as soon as possible, it then takes the first available el. connection to the next city v_{i+1} . Since the el. connections do not overtake, Expand reached the city v_{i+1} as soon as possible.

□

We would like to stress that overtaking is understood as a situation when one carrier overtakes another between *two subsequent stations*. This situation is not that common, however it is still present in the real world timetables ¹³, as shown in table 5.1. All the same, we can simply remove the

¹³In Slovak rails, no overtaking has been detected. This is not surprising as (to my knowledge) there are no inter-station tracks with multiple rails going in one direction. French railways, on the other hand have designated high-speed tracks and thus overtaking is not impossible.

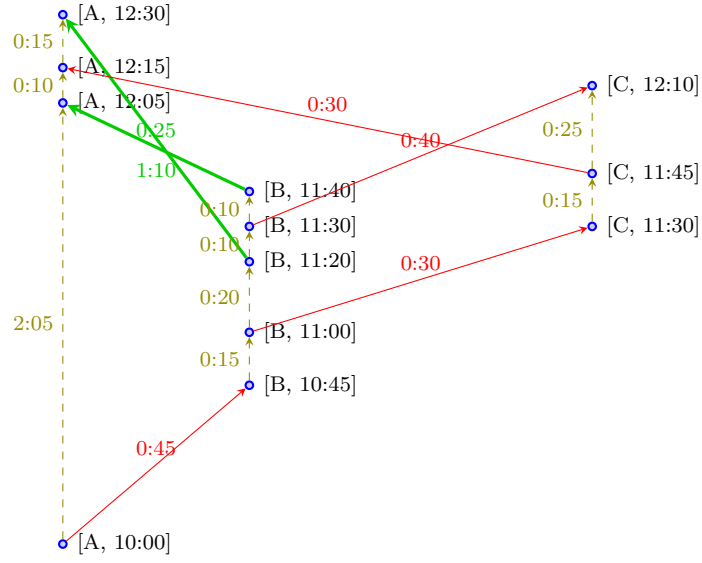


Figure 5.2: An example of **overtaking** (in thick), depicted in a TE graph.

overtaken el. connections from the timetables, as they can be substituted by the quicker connection plus some waiting.

Name	Overtaken edges (%)
<i>air01</i>	1%
<i>cpsk</i>	2%
<i>gb-coach</i>	1%
<i>gb-train</i>	0%
<i>montr</i>	1%
<i>sncf</i>	1%
<i>sncf-ter</i>	1%
<i>sncf-inter</i>	6%
<i>zsr</i>	0%

Table 5.1: Presence of overtaking in the timetables.

The basic idea of the algorithm *USP-OR* (a short-cut for USP oracle) is therefore simply to pre-compute all the USPs for each pair of cities. Upon a query, the algorithm simply expands all the USPs for a given pair of cities, reconstructs respective connections and chooses the best one.

Algorithm 5.2 *USP-OR* query

Input

- timetable T
- OC query (x, t, y)

Pre-computed

- $\forall x, y$: set of USPs between x and y ($usps(x, y)$)

Algorithm

```
 $c^* = null$ 
for all  $p \in usps_{x,y}$  do
   $c = Expand(T, p, t)$ 
   $c^* = \text{better out of } c^* \text{ and } c$ 
end for
```

Output

- connection c
-

5.1.1 Analysis of USP-OR

We will now have a look at the four parameters of this oracle based method. As for the preprocessing time, we need to find optimal connections from each *event* in the timetable to each *city* (or in other words - solve all possible OC queries). On these connections we apply the *path* function to obtain the USPs. The maximum number of events in one city is the height h and there is n cities, thus hn is the upper bound on the number of events. One search from a single event to all cities can be done in time $\mathcal{O}(n \log n + m)$ with a TD Dijkstra's algorithm run on the time-dependent graph of our timetable (TD_T). In worst case, m could be as much as n^2 but we may bound it as $m \leq \delta_T n$ (where δ_T is the sparsity of the timetable, defined in section 4). We therefore get the **preprocessing time** $\mathcal{O}(hn^2(\log n + \delta))$.

As for the preprocessed space, we need to store USPs for each pair of the cities (n^2 pairs) and each USP might be long at most $\mathcal{O}(n)$ hops. What is more, there might be many USPs for a single pair of cities. Therefore we have two questions with respect to the space complexity of the preprocessing:

1. What is the average size of the USPs?
2. How many are there USPs between a single pair of cities?

The answer for the first question is that the average USP size is equal to the OC radius of the timetable (γ_T) defined in section 4, which generally varies around \sqrt{n} (see table ??).

To answer the second question, we will introduce the following definition:

Definition 5.3. USP coefficient

Given a timetable T and a pair of cities x, y , the USP coefficient $\tau_T(x, y) = |usps_T(x, y)|$, where $usps_T(x, y)$ is the set of USPs between x and y . By τ_T we will denote the average USP coefficient in timetable T .

From the table 5.2 we can see, that there are not many USPs on average, meaning that τ is usually some small number. Also, we see that it slightly increases with increasing time range (plot 5.6), but not with increasing n , the size of the timetable (plot 5.3). Thus we can consider τ to be bound by a small constant when it comes to daily timetables.

From the answers to our two questions we see that the **size of the preprocessed oracle** is $\mathcal{O}(\tau n^2 \gamma)$.

Name	τ	$\max \tau(x, y)$
<i>air01-d</i>		
<i>cpsk-d</i>		
<i>gb-coach-d</i>		
<i>gb-train-d</i>		
<i>montr-d</i>		
<i>sncf-d</i>		
<i>sncf-ter-d</i>		
<i>sncf-inter-d</i>		
<i>zsr-d</i>		

Table 5.2: Average and maximal USP coefficients for daily timetables.

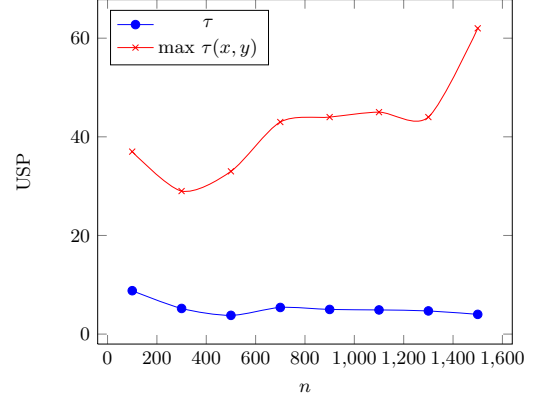


Figure 5.3: Changing of τ with increased number of stations in *sncf* dataset.

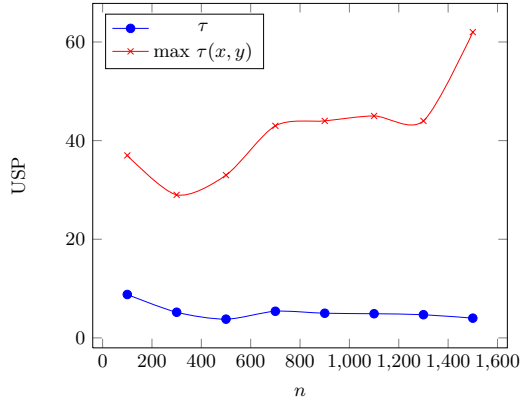


Figure 5.4: Changing of τ with increased number of stations in *cpsk* dataset.

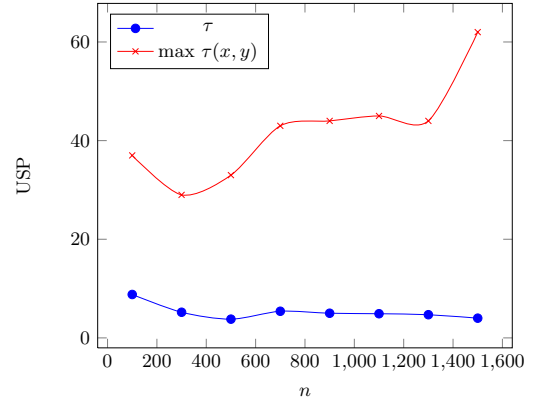


Figure 5.5: Changing of τ with increased number of stations in *gb-coach* dataset.

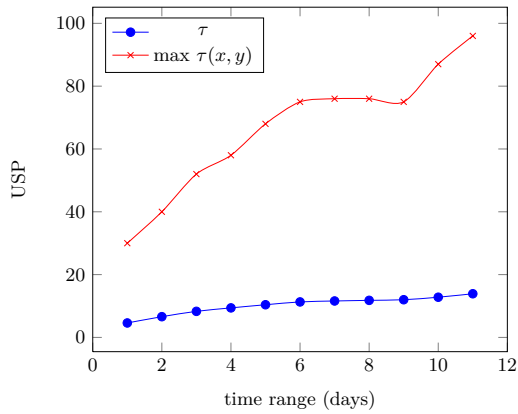


Figure 5.6: Changing of τ with increased time range in *air01* dataset. 1 day = about 800 in height.

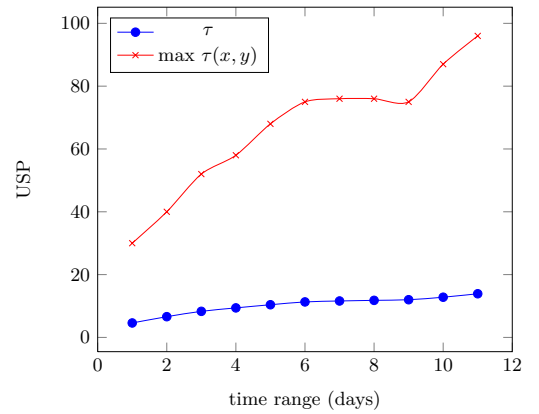


Figure 5.7: Changing of τ with increased time range in *zsr* dataset. 1 day = about in height.

The query time also depends on the USP coefficient of a given pair of cities x, y , as we have to try out all USPs in $usps(x, y)$. The expansion of a USP by *Expand* function takes time linear in the size of the USP ¹⁴, leading to **query time** $\mathcal{O}(\tau\gamma)$ on average. Note, that this is pretty much optimal, as τ is basically constant and we need to output the connection itself, which takes linear time in its size.

Finally, the **stretch** of *USP-OR* is **1**, as it returns exact answers.

<i>USP-OR</i>	<i>prep</i>	<i>size</i>	<i>qtime</i>	<i>stretch</i>
guaranteed	$\mathcal{O}(hn^2(\log n + \delta))$	$\mathcal{O}(\tau n^2 \gamma)$	avg. $\mathcal{O}(\tau\gamma)$	1
τ const., $\gamma \leq \sqrt{n}$, $\delta \leq \log n$	$\mathcal{O}(hn^2 \log n)$	$\mathcal{O}(n^{2.5})$	avg. $\mathcal{O}(\sqrt{n})$	1

Table 5.3: The summary of the *USP-OR* algorithm parameters. The second row corresponds e.g. to the *sncl* dataset.

5.2 USP-OR-A

With *USP-OR* the main disadvantage is its space consumption. We may decrease this space complexity by pre-computing USPs only among *some* cities. The nodes that we select for this purpose will be called **access nodes** (AN for short), as for each city they would be the crucial nodes we need to pass in order to access most of the cities of T . It would be suitable for this access node set to have several desirable properties. In order to formulate them, we need to define a few terms first.

Definition 5.4. Front neighbourhood

Given a timetable T and access node set \mathcal{A} , a *front neighbourhood* of city x are all cities (including x) that are reachable from x not via \mathcal{A} . Formally $\mathbf{neigh}_{\mathcal{A}}(x) = \{y \mid \exists \text{ path } p = (p_1, p_2, \dots, p_k) \text{ from } x \text{ to } y \text{ in } ug_T : p_i \neq a \forall a \in \mathcal{A}, i \in \{2, \dots, k-1\}\}$ ¹⁵

We define analogically **back neighbourhood** (denoted $\mathbf{bneigh}_{\mathcal{A}}(x)$), as nodes that could be reached in reversed UG ($\overleftarrow{ug_T}$). Note that the access nodes that are on the boundary of x 's neighbourhoods are also part of these neighbourhoods. These access nodes form some sort of separator between the x 's neighbourhood and the rest of the graph and we will call them **local access nodes (LAN)** ($\mathbf{lan}_{\mathcal{A}}(x) = \mathcal{A} \cap \mathbf{neigh}_{\mathcal{A}}(x)$), or analogically **back local access nodes (blan_A(x))**.

Now we may formulate the three desired properties of the access node set \mathcal{A} . Given a timetable T and small constants r_1, r_2 and r_3 , we would like to find access node set \mathcal{A} such that:

1. The access node set is sufficiently small

$$|\mathcal{A}| \leq r_1 \cdot \sqrt{n} \quad (5.1)$$

2. The average square of neighbourhood ¹⁶ size for cities not in \mathcal{A} is at most $r_2 \cdot n$

$$\frac{\sum_{x \in ct_T \setminus \mathcal{A}} |\mathbf{neigh}_{\mathcal{A}}(x)|^2}{|ct_T \setminus \mathcal{A}|} \leq r_2 \cdot n \quad (5.2)$$

¹⁴In time-dependent graphs, this requires a constant-time retrieval of the correct interpolation point of the cost function (the piece-wise linear function that tells us the traversal time of an arc at a given time) for some time t . More specifically, we need to obtain an interpolation point $\mathit{argmin}_{(t', l)} \{t' \mid t' > t\}$. If we assume uniform distribution of departures throughout the time range of the timetable, this can be implemented in constant time. Otherwise, binary search lookup is possible in time $\mathcal{O}(\log h)$.

¹⁵We leave out subscript identifying the timetable T . In situation with clear context, we may also leave out the \mathcal{A} subscript.

¹⁶We required the same for back neighbourhoods.

3. The average square of the number of local access nodes ¹⁷ for cities not in \mathcal{A} is at most r_3

$$\frac{\sum_{x \in ct_T \setminus \mathcal{A}} |lan_{\mathcal{A}}(x)|^2}{|ct_T \setminus \mathcal{A}|} \leq r_3 \quad (5.3)$$

An access node set \mathcal{A} with the above mentioned properties will be called **(r_1, r_2, r_3) access node set** (AN set). We will now explain how the *USP-OR-A* (*USP-OR* with access nodes) algorithm works and return to its analysis later.

During preprocessing, we need to find a good AN set and compute the USPs between every pair of access nodes. For every city $x \notin \mathcal{A}$, we also store its $neigh_{\mathcal{A}}(x)$, $bneigh_{\mathcal{A}}(x)$, $lan_{\mathcal{A}}(x)$ and $blan_{\mathcal{A}}(x)$. On a query from x to y at time t , we will first make a local search in the neighbourhood of x up to x 's local access nodes. Subsequently, we want to find out the earliest arrival times to each of y 's *back* local access nodes. To do this, we take advantage of the pre-computed USPs between access nodes - try out all the pairs $u \in lan(x)$ and $v \in blan(y)$ and expand the stored USPs. Finally, we make a local search from each of y 's back LANs to y , but we run the search *restricted* to y 's back neighbourhood. For more details, see algorithms 5.3 and 5.4 and picture 5.8, where we have split the algorithms to 3 distinct phases.

Algorithm 5.3 *USP-OR-A* preprocessing

Input

- timetable T

Algorithm

find a good AN set \mathcal{A}

$\forall x, y \in \mathcal{A}$ compute $usps(x, y)$

$\forall x \in ct_T \setminus \mathcal{A}$ compute $neigh_{\mathcal{A}}(x)$, $bneigh_{\mathcal{A}}(x)$, $lan_{\mathcal{A}}(x)$ and $blan_{\mathcal{A}}(x)$

Output

- output everything we have computed
-

¹⁷We required the same for back LANs.

Algorithm 5.4 *USP-OR-A* query

Input

- timetable T
- OC query (x, t, y)

Algorithm

let $lan(x) = x$ if $x \in \mathcal{A}$

let $blan(y) = y$ if $y \in \mathcal{A}$

Local front search

perform TD Dijkstra from x at time t up to $lan(x)$

if $y \in neigh(x)$ **then**

 let c_{loc}^* be the connection to y obtained by TD Dijkstra *# the optimal connection may still go via ANs (though it is unlikely)*

end if

$\forall u \in lan(x)$ let $ea(u)$ be the arrival time and $oc(u)$ the conn. to u obtained by TD Dijkstra

Inter-AN search

for all $v \in blan(y)$ **do**

$oc(v) = null$

for all $u \in lan(x)$ **do**

for all $p \in usps(u, v)$ **do**

$c = Expand(T, p, ea(u))$

$oc(v) = \text{better out of } oc(v) \text{ and } c$

end for

end for

end for

$\forall v \in blan(y)$ let $ea(v) = end(oc(v))$

Local back search

for all $v \in blan(y)$ **do**

 perform TD Dijkstra from v at time $ea(v)$ to y restricted to $bneigh(y)$

 let $fin(v)$ be the connection returned by TD Dijkstra

end for

$v^* = \operatorname{argmin}_{v \in blan(y)} \{end(fin(v))\}$

$u^* = from(oc(v^*))$

let $c^* = oc(u^*).oc(v^*).fin(v^*)$ *# the dot (.) symbol is concatenation of connections*

output better out of c_{loc}^* and c^*

Output

- optimal connection $c_{(x,t,y)}^*$
-

5.2.1 Analysis of USP-OR-A

Let us now analyse the properties of this oracle-based method. Clearly, much depends on the way we look for the access node set. We will address this issue in next subsections but for now, we will assume we can find (r_1, r_2, r_3) AN set \mathcal{A} in time $f(n)$. Then, in the preprocessing, we have to find USPs among the access nodes, which requires running Dijkstra's algorithm from each event in a city from \mathcal{A} . There is $\mathcal{O}(r_1 h \sqrt{n})$ such events which leads to the time complexity $\mathcal{O}(r_1 h n^{1.5} (\log n + \delta))$. We also have to find local access nodes and neighbourhoods for each city, which can be accomplished with e.g. depth first search exploring the neighbourhood. This search algorithm (run from non-access city) has complexity linear in the number of arcs and so we could bound the total complexity as:

$$\sum_{x \in ct_T \setminus \mathcal{A}} |E(neigh_{\mathcal{A}}(x))| \leq \sum_{x \in ct_T \setminus \mathcal{A}} |neigh_{\mathcal{A}}(x)|^2 \leq r_2 n^2$$

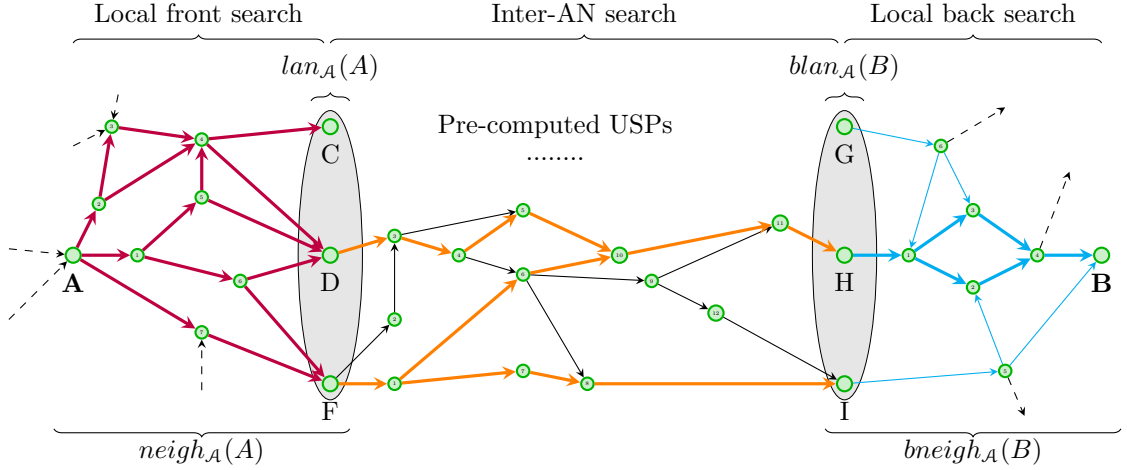


Figure 5.8: Principle of *USP-OR-A* algorithm. The arcs in **bold** mark areas that will be explored: all nodes in $neigh_A(x)$, USPs between LANs of x and back LANs of y and the back neighbourhood of y (possibly only part of it will be explored, since the local back search goes against the direction in which the back neighbourhood was created).

where $E(V)$ is the set of arcs among vertices of V . However this is very loose upper bound, as our UGs are actually very sparse. Therefore we can improve it. We know from the equation 5.2 that the average square of neighbourhood size is $\leq r_2 \cdot n$. As a consequence of the Cauchy-Schwarz Inequality [?] the following holds for positive real numbers x_i :

$$\sqrt{\frac{x_1^2 + x_2^2 + \dots + x_n^2}{n}} \geq \frac{x_1 + x_2 + \dots + x_n}{n}$$

Applying this to our neighbourhood sizes, we get that the average size of the neighbourhood is at most $\sqrt{r_2 n}$. We now split the vertices of $ct_T \setminus \mathcal{A}$ to two categories: those with neighbourhoods of size $\leq \sqrt[4]{n}$ will be part of the set S_{\leq} and those with neighbourhoods of size bigger then $\sqrt[4]{n}$ will be in $S_{>}$. A neighbourhood in the first category cannot possibly contain more than \sqrt{n} arcs while those in the second category can have at most $\delta_T |neigh_A(x)|$ arcs, depending on the timetable's density.

$$\begin{aligned} \sum_{x \in ct_T \setminus \mathcal{A}} |E(neigh_A(x))| &\leq \\ \sum_{x \in S_{\leq}} \overbrace{|E(neigh_A(x))|}^{\leq \sqrt{n}} + \sum_{x \in S_{>}} \overbrace{|E(neigh_A(x))|}^{\leq \delta |neigh_A(x)|} &\leq \\ n\sqrt{n} + \delta n\sqrt{r_2 n} &\leq \\ \delta r_2 n^{1.5} \end{aligned}$$

Therefore, the total **time complexity of the preprocessing** is $\mathcal{O}(f(n) + r_1 h n^{1.5} (\log n + \delta)) + \mathcal{O}(\delta r_2 n^{1.5}) = \mathcal{O}(f(n) + (r_1 + r_2)(\delta + \log n) h n^{1.5})$.

As for the size of the preprocessed data - we need to store all the neighbourhoods, LANs and USPs between access nodes. We already know that the average size of the neighbourhood is $\leq \sqrt{r_2 n}$, thus the total size of the (front and back) neighbourhoods is $\mathcal{O}(r_2 n^{1.5})$ ¹⁸. This term

¹⁸As r_2 will be a very small constant, we may disregard the square root.

bounds also the size of the pre-computed local access nodes for each node.

Finally we have the preprocessed USPs. There is at most $r_1^2 n$ pairs of access nodes and for each of them we have possibly several USPs. We will denote by $\tau_{\mathcal{A}}$ the average USP coefficient between pairs of cities from \mathcal{A} and by $\gamma_{\mathcal{A}}$ the average optimal connection size (or equivalently, USP size) between cities in \mathcal{A} . This amounts to $\mathcal{O}(r_1^2 \tau_{\mathcal{A}} \gamma_{\mathcal{A}} n)$ for storage of USPs and to a total **preprocessing size** $\mathcal{O}(r_2 n^{1.5} + r_1^2 \tau_{\mathcal{A}} \gamma_{\mathcal{A}} n)$.

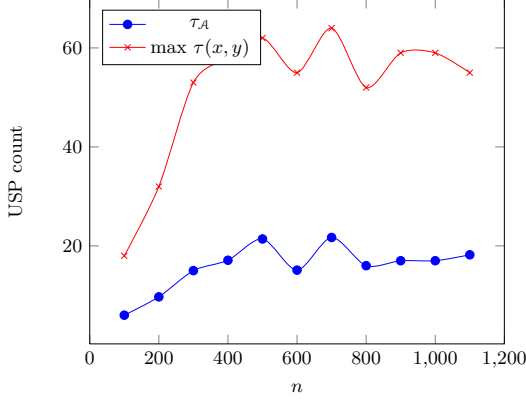


Figure 5.9: Changing of $\tau_{\mathcal{A}}$ with increased number of stations in *cpsk* dataset. \mathcal{A} was obtained using algorithm *Locsep* we will talk about later.

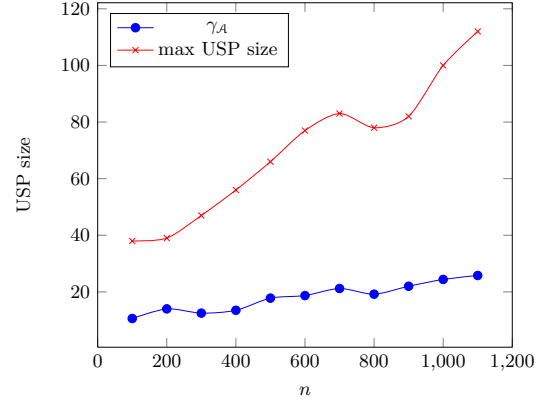


Figure 5.10: Changing of $\gamma_{\mathcal{A}}$ with increased number of stations in *cpsk* dataset. \mathcal{A} was obtained using algorithm *Locsep* we will talk about later.

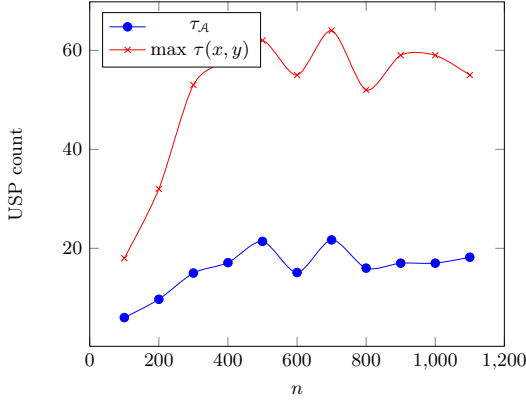


Figure 5.11: Changing of $\tau_{\mathcal{A}}$ with increased number of stations in *gb-train* dataset. \mathcal{A} was obtained using algorithm *Locsep* we will talk about later.

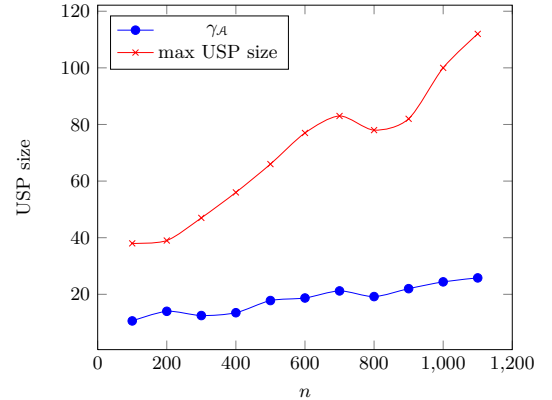


Figure 5.12: Changing of $\gamma_{\mathcal{A}}$ with increased number of stations in *gb-train* dataset. \mathcal{A} was obtained using algorithm *Locsep* we will talk about later.

On a query from x at time t to y , we first perform the *local front search* (see algorithm 5.4). In this step we explore the neighbourhood of x with a time-dependent Dijkstra's algorithm, which takes on average time $\mathcal{O}(\sqrt{r_2 n}(\log(\sqrt{r_2 n}) + \delta))$. We then expand all the USPs between u and v such that $u \in \text{lan}(x)$ and $v \in \text{blan}(y)$, which takes on average $\mathcal{O}(r_3 \tau_{\mathcal{A}} \gamma_{\mathcal{A}})$. Finally, from each $v \in \text{blan}(y)$ we do a TD Dijkstra, restricted to $\text{bneigh}(y)$, leading to time complexity $\mathcal{O}(r_3 \sqrt{r_2 n}(\log(\sqrt{r_2 n}) + \delta))$.

Summing up the three terms we obtain the **query time** of $\mathcal{O}(r_2 r_3 \sqrt{n}(\log(r_2 n) + \delta) +$

$r_3\tau_A\gamma_A$).

Stretch of the *USP-OR-A* algorithm is **1**, as it is exact algorithm.

The resulting bounds do not look very appealing. This is because we wanted to preserve the generality - the concrete bounds will depend on what kind of properties the timetables have and what algorithm for finding the AN set is plugged in. In table 5.4, we summarize the parameters of *USP-OR-A* method and provide the bounds for a case when the properties of the timetables correspond to those we have measured in our datasets and when we have an algorithm that finds good AN set.

<i>USP-OR-A</i>	guaranteed	τ, r_1, r_2, r_3 const., $\gamma \leq \sqrt{n}, \delta \leq \log n$
<i>prep</i>	$\mathcal{O}(f(n) + (r_1 + r_2)(\delta + \log n)hn^{1.5})$	$\mathcal{O}(f(n) + hn^{1.5} \log n)$
<i>size</i>	$\mathcal{O}(r_2n^{1.5} + r_1^2\tau_A\gamma_An)$	$\mathcal{O}(n^{1.5})$
<i>qtime</i>	avg. $\mathcal{O}(r_2r_3\sqrt{n}(\log(r_2n) + \delta) + r_3\tau_A\gamma_A)$	avg. $\mathcal{O}(\sqrt{n} \log n)$
<i>stretch</i>	1	1

Table 5.4: The summary of the *USP-OR-A* algorithm parameters.

5.2.2 Correctness of *USP-OR-A*

Finally, we will proof the correctness of the algorithm, i.e. that it always returns the optimal connection.

Theorem 5.1. *The algorithm USP-OR-A 5.3 5.4 always returns the optimal connection.*

Proof. Let \mathcal{A} be the set of access nodes and consider a query from city x to city y at any time t . If $x \in \mathcal{A}$ and $y \in \mathcal{A}$, an optimum is returned due to lemma 5.1 (in such a case, we basically run *USP-OR* algorithm).

In the following we will assume that $y \notin \text{neigh}(x)$, which means that the optimal connection goes through some access node $u \in \text{lan}(x)$ and $v \in \text{blan}(y)$. Note that it may be that $u = v$.

What we would like to prove as a next step is that we reach the back LANs of y (or y itself if it is an access node) at the earliest arrival time. After the *local front search*, we have reached the x 's local ANs at times $ea(u) \forall u \in \text{lan}(x)$. For some local access node this value is the true earliest arrival. Let us denote the set of such local ANs as $\text{lan}^*(x)$. The crucial thing to realize is, that the optimal connection to any city out of the x 's neighbourhood will lead via some $u \in \text{lan}^*(x)$ (see picture 5.13). And because the *inter-AN search* phase finds *optimal* connections between pairs $u \in \text{lan}(x)$ and $v \in \text{blan}(y)$, it follows that for each $v \in \text{blan}(y)$ the $ea(v)$ is the earliest arrival to this city after the *inter-AN search* phase.

In the *local back search* we run a TD Dijkstra search from all back LANs of y . And since this algorithm is exact and starts from each back LAN as early as possible, we get the optimal connection to y .

It remains to show that if $y \in \text{neigh}(x)$, we also get the optimal connection. In such case, we simply compare the connection that goes via access nodes and the one that was obtained solely within the neighbourhood and output the shorter one. As there are no other options, the proof is complete. \square

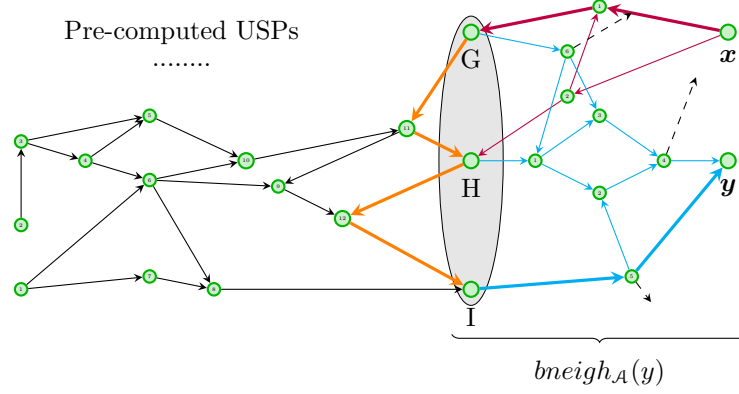


Figure 5.13: On the picture $lan(x) = \{G, H\}$ and $blan(y) = \{G, H, I\}$. In **thick** we have highlighted the optimal connection. The connection to H is sub-optimal after the *local front search* phase, however the optimal connection to y (and to H and I as well) leads through $lan^*(x)$ (some of x 's local access nodes to which we have an optimal connection after the *local front search*. Particularly, it goes through G).

5.2.3 Modifications of *USP-OR-A*

Our implementation of the *USP-OR-A* algorithm uses one slight improvement, which we did not mention in its description, since it is more of a optimization technique without any theoretical guarantees on actual improvement of the running time. However, we consider it an interesting idea so we mention it at this place.

Definition 5.5. *USP tree*

Given a pair of cities x and y in a timetable T , we will call a *USP tree* the graph made out of edges of all USPs in $usp_T(x, y)$: $usp_T^3(x, y) = (V^3, E^3)$ where $V^3 = \{v \mid v \text{ lays on some } p \in usp_T(x, y)\}$ and $E^3 = \{(a, b) \mid (a, b) \text{ is part of some } p \in usp_T(x, y)\}$.

We could take advantage of these USP trees to speed up the *local front search* phase of the algorithm, where we unnecessarily explore the whole neighbourhood when we could just go along the arcs of the USP trees. The picture 5.14 depicts this.

The interesting thing about this is the exploitation of both - timetable and its underlying graph. While the neighbourhood of a node is something static, related only to the structure of the UG and generally time-independent, the USP trees reflect to some extent the properties of the timetable (e.g. which ways are frequently serviced and thus provide optimal connections). By intersecting these two things, we get the area that is *worth* to be explored and that is *small* at the same time (provided, of course, that the neighbourhoods are small).

5.3 Selection of access node set

The challenge in the *USP-OR-A* algorithm comes down to the selection of a good access node set - a (r_1, r_2, r_3) AN set with both three parameters as low as possible. However, intuitively (and experimentally verified), decreasing e.g. r_1 (the AN set size) increases r_2 (the size of the neighbourhoods). We therefore have to do some compromises.

In the following we first show the problem of choosing an optimal access node set to be NP-hard. We then present our methods for heuristic selection of access nodes and show their performance on real data.

all vertices of m_j to $s_i \iff j \in S_i$. Finally, for we connect s_i to s'_i , $1 \leq i \leq k$.

Example. Let $m = 10$ (thus $U = \{1, 2, \dots, 10\}$) and $k = 13$:

- $S_1 = \{1, 3, 10\}$
- $S_2 = \{1, 2\}$
- ...
- $S_{13} = \{2, 3, 10\}$

For this instance of min set-cover, we construct the graph depicted on picture 5.15.

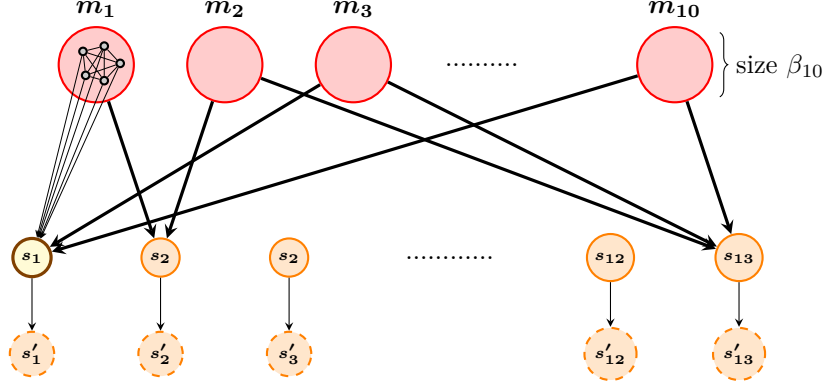


Figure 5.15: The principle of the reduction. In m_i , there are actually complete graphs of β_i vertices (as shown for m_1). **Thick** arcs represent arcs from all the vertices of respective m_i . The s_i vertices are connected to their s'_i versions. If e.g. s_1 is selected as an access node, s'_1 is no longer part of any neighbourhood.

Now we would like to clarify the sizes of m_i . Define α_i to be the number of sets S_j that contain i : $\alpha_i = |\{S_j \in \mathcal{S} \mid i \in S_j\}|$ and assume the constructed graph has n vertices. We want the β_i to satisfy $\beta_i \geq 2$ and $\beta_i + 2\alpha_i - 1 \leq \sqrt{n}$ but $\beta_i + 2\alpha_i > \sqrt{n}$. The last two inequalities would mean that if at least one s_j connected to m_i is chosen as an access node, the neighbourhood for nodes in m_i will be still $\leq \sqrt{n}$, but if none of them is chosen, the neighbourhood will be just over \sqrt{n} . For now we will assume that we have constructed the graph in such a way that all β_i satisfy the mentioned inequalities. We will return to construction of the graph at the end.

Now consider an optimal AN set which contains a vertex from within some m_i . If this is the case, **either** some s_j to which m_i is connected is selected as AN, **or** all vertices from m_i are access nodes **or** the neighbourhood is too large. Keep in mind that the local access nodes are also part of neighbourhoods, so unless we select for AN some of the s_j that m_i is connected to, the neighbourhood of any non-access node in m_i will be too large. As there are at least two nodes in every m_i , it is more efficient to select some s_j rather than select all nodes in m_i . Thus when it comes to selecting ANs *it is worth to consider only vertices s_j* .

From this point on, it is easy to see that it is optimal to select those s_j that correspond to the optimal solution of min-set cover. The reason is that each of the m_i will be connected to at least one access node s_j and will thus have neighbourhood size $\leq \sqrt{n}$, while the number of selected access nodes will be optimal.

It remains to show how to choose values β_i . Due to the condition $\beta_i \leq \sqrt{n} - 2\alpha_i + 1$ we need to have sufficiently big n to fulfil $\beta_i \geq 1$. We will accomplish this by adding dummy isolated vertices

to the graph. Define function $nextSquare(x)$ to output the smallest $y^2 > x$ where y is a natural number. We then compute $w = (\max\{2\alpha_i\} + 2)^2$ and select the starting value of n to be $n' = nextSquare(\max\{w - 1, \sqrt{2k + m}\})$. We create the s_j and s'_j vertices and complete graphs m_i containing so far only one vertex each. We connect everything according to the rules stated earlier in this proof and we create dummy vertices up to the capacity defined by n . Now we repeat the following:

- We compute \sqrt{n} which is a natural number
- For i from 1 to m we add vertices to m_i till it does not contain $\sqrt{n} - 2\alpha_i + 1$ vertices. For each added vertex we delete one dummy vertex.
- If we run out of dummy vertices, $n = nextSquare(n)$
- Break out of the loop if $|m_i| = \sqrt{n} - 2\alpha_i + 1 \forall i$

With each iteration of this little algorithm we will be forced to add one more vertex to all m_i (since \sqrt{n} increased by one), a so called *inefficient increase*. At the beginning, we need to make at most $m\sqrt{n'}$ efficient increases to meet the breaking condition. And since m is constant and the capacity of new dummy vertices increases linearly, after t steps we create $\mathcal{O}(t^2)$ dummy vertices that may be used for efficient increases. Therefore, the algorithm will stop after $\mathcal{O}(\sqrt{mn'})$ steps. \square

5.3.2 Choosing ANs based on node properties

In the previous sub-subsection, we have shown the problem of choosing the optimal AN set to be NP-hard. In this sub-subsection we perform a simple experiment of choosing for the access nodes the cities that seem to be the most important. More specifically, in the optimistic underlying graph (see section 2) ug_T^{opt} we were looking for cities with:

1. High **degree**. We consider the sum of in-degree and out-degree ¹⁹ of the respective node x :
 $deg(x) = deg_{in}(x) + deg_{out}(x)$.
2. High **betweenness centrality** (BC). Betweenness centrality for a node v is defined as

$$g(v) = \sum_{s \neq v \neq t} \frac{\sigma_{st}(v)}{\sigma_{st}}$$

where $\sigma_{st}(v)$ is the number of shortest paths from s to t passing through v and σ_{st} is the total number of shortest paths from s to t [?]. We then scale the values to the range $< 0, 1 >$ to obtain for each city x its scaled betweenness centrality $bc(x)$.

We will denote by $\mathcal{A}_{deg}(k)$ the set of k cities with highest $deg(x)$ value. We were interested in the smallest k such that:

1. $\mathcal{A}_{deg}(k)$ is (r_1, r_2, r_3) AN set with $r_2 \leq 1$ (the average square of neighbourhoods is $\leq \sqrt{n}$). Denote such r_1 as \mathbf{r}_{deg}^{avg} ²⁰.
2. $\mathcal{A}_{deg}(k)$ is (r_1, r_2, r_3) AN set where $\forall x \in ct_T: |neigh_{\mathcal{A}_{deg}(k)}(x)| \leq \frac{3\sqrt{n}}{2}$ and $|bneigh_{\mathcal{A}_{deg}(k)}(x)| \leq \frac{3\sqrt{n}}{2}$ (all neighbourhoods are large at most $\frac{3\sqrt{n}}{2}$). Denote such r_1 as \mathbf{r}_{deg}^{max} .

We define similarly \mathbf{r}_{bc}^{avg} and \mathbf{r}_{bc}^{max} . The resulting values for datasets *sncf* and *cpru* could be seen on plots 5.16.

¹⁹In-degree is the number of arcs going towards to node and out-degree the number of outgoing arcs.

²⁰Intuitively - \mathbf{r}_{deg}^{avg} is the smallest r_1 such that $r_1\sqrt{n}$ highest-degree nodes selected as ANs are enough to satisfy that the average square of neighbourhoods is $\leq \sqrt{n}$.

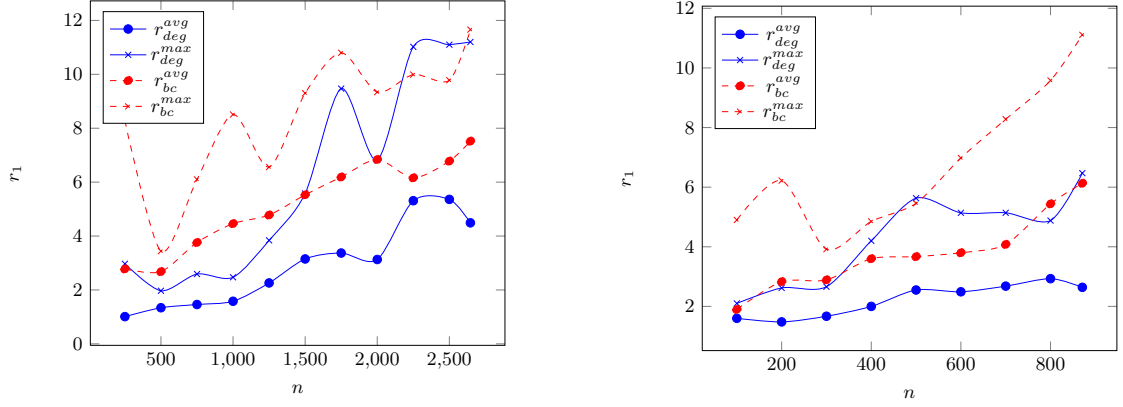


Figure 5.16: Necessary sizes of access node sets ($|\mathcal{A}| = r_1\sqrt{n}$) based on high-degree/high-BC cities. Datasets *snkf* (left) and *cpsk* (right). Notice the occasional “roller coaster” bumps (especially on the left plot) - an explanation of this phenomena is that in the immediately smaller sub-timetable we have erased just the high-degree node that proved to be a good access node, and which now must be substituted by many other access nodes.

5.3.3 Choosing ANs heuristically - the *Locsep* algorithm

Clearly, selecting the cities for access nodes solely by high degree or BC value is not the best way. Probably the few nodes with highest degrees and BC will indeed be part of the AN set, as they are intuitively some sort of central hubs without which the network would not work. However, after we select these most important nodes to the AN set, we need some better measure of node’s importance, or suitability to be an access node. In the following we present a simple heuristic approach run on underlying graph ug_T of given timetable T that evaluates its vertices based on how good local separators they are.

The algorithm that we call *Locsep* (as it looks for good local separators) will work in iterations, each of them resulting in a selection of the city with the highest score to the access node set \mathcal{A} ²¹. We continue to select access nodes until we meet the following stopping criterion: \mathcal{A} is (r_1, r_2, r_3) AN set with $r_2 \leq 1$ (the average square of neighbourhoods is $\leq \sqrt{n}$) ²². We will denote $r_{ls} = r_1$ of the resulting access node set.

The important thing that remains to be shown is how do we compute the score for a particular city. The following text explains this.

In each iteration, we first compute the neighbourhoods and back neighbourhoods (given the current access node set \mathcal{A}) for each city. We need this to evaluate the stopping criteria, but the information is also used in the computation of the **potential** (the score) of the cities.

For a city x , we compute its potential p_x in the following way: we explore an area A_x of \sqrt{n} nearest cities around x , ignoring branches of the search that start with an access node (x is an exception to this, since we start the search from it. However $x \notin A_x$ holds). We do this exploration in an underlying graph with no orientation and no weights. Next we get the front and back neighbourhoods of x within A_x ($\mathbf{fn}(x) = \text{neigh}(x) \cap A_x$, $\mathbf{bn}(x) = \text{bneigh}(x) \cap A_x$).

²¹Actually, in our implementation, we allow an occasional deselection of an already selected node with the *lowest* score, to avoid having in the resulting set cities that had high score when selected but were not very useful access nodes at the end.

²²In our implementation, we perform some further adjustments of the resulting set, such as removing unnecessary access nodes and optimising the r_3 value.

For a set of access nodes \mathcal{A} , let us call a path p in ug_T **access-free** if it does not contain a node from \mathcal{A} . Now as long as x is not in \mathcal{A} , we have a guarantee that for every pair $u \in bn(x)$ and $v \in fn(x)$ there is an access-free path from u to v within A_x . Our interest is how this will change after the selection of x .

Consider now a node $y \in bn(x)$. We will call $\mathbf{sur}(y) = \max\{0, |neigh(y)| - \sqrt{n}\}$ the **surplus** of y 's neighbourhood, i.e., by how much we wish to reduce it so that it is $\leq \sqrt{n}$. If the surplus is zero, y will not add anything to the x 's potential. Otherwise, we run a restricted (to A_x) search from y during which we explore j vertices in $fn(x)$. We increase the potential of x by $\min\{\mathbf{sur}(y), |fn(x) - j|\}$ - i.e. by how much we can decrease the surplus of y 's neighbourhood. We do the same for all $y \in bn(x)$ and a similar thing for all $y \in fn(x)$ (we use $\overleftarrow{ug_T}$ instead of ug_T , $bneigh(y)$ instead of $neigh(y)$ etc...). For an illustration of potential computing, see picture 5.17.

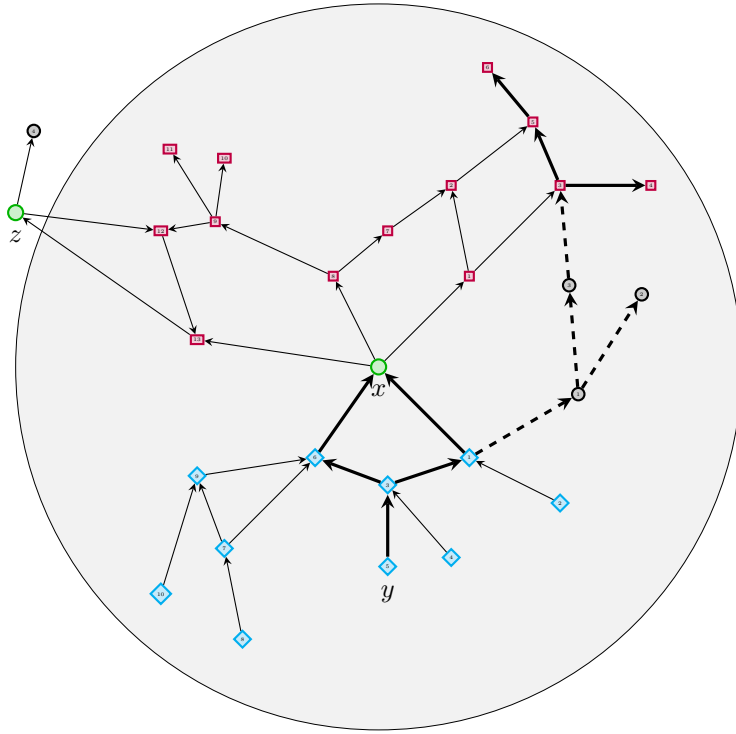


Figure 5.17: The principle of computing potentials in *Locsep* algorithm. We explored an area of \sqrt{n} nearest cities (in terms of hops) around x . Access nodes (like z) and cities behind them are ignored. Little **squares** are nodes from $fn(x)$ and **diamonds** are part of $bn(x)$. From y we run a forward search (the **thick** arcs). Nodes from the $fn(x)$ that were not explored in this search can only be reached via x itself. Such nodes contribute to x 's potential assuming y has large neighbourhood size.

Finally, we simply get the city $x \notin \mathcal{A}$ with the highest potential and select it as an access node. We check the stopping criterion and in case it is not satisfied yet, we move on to next iteration. However, note that when a new node x' is selected to \mathcal{A} , we do not have to re-compute neighbourhoods and potentials of all cities - it is only necessary for those cities that could reach/be reached access-free from x' (i.e. nodes from $neigh_{\mathcal{A}}(x') \cup bneigh_{\mathcal{A}}(x')$). Algorithm 5.5 provides a high-level overview of the *Locsep* method.

Algorithm 5.5 *Locsep*

Input

- ug_T

Algorithm $\mathcal{A} = \emptyset$ $ct' = ct_T$ **while** $r_2 > 1$ **do** $\forall x \in ct'$: compute $neigh_{\mathcal{A}}(x)$, $bneigh_{\mathcal{A}}(x)$ $\forall x \in ct'$: compute p_x $x' = \operatorname{argmax}_{x \notin \mathcal{A}} \{p_x\}$ $\mathcal{A} = \mathcal{A} \cup \{x'\}$ $ct' = neigh_{\mathcal{A}}(x') \cup bneigh_{\mathcal{A}}(x')$ **end while**

Remove unnecessary ANs

Optimise r_3 **Output**

- AN set \mathcal{A} , such that $|\mathcal{A}| = \sqrt{nr_{locsep}}$
-

Now we would like to estimate the **time complexity** of Locsep algorithm. As mentioned, one iteration consists of three parts:

1. Computing neighbourhoods. Unfortunately, at the beginning when $\mathcal{A} = \emptyset$, the neighbourhood sizes may be as large as $\mathcal{O}(n)$. Therefore, we may bound the complexity of this phase only as $\mathcal{O}(nm) = \mathcal{O}(\delta n^2)$.
2. Computing potentials. For a city x we explore area of the size \sqrt{n} and from each node in that area we do a restricted search. Therefore the total complexity of this step is $\mathcal{O}(n \cdot \sqrt{n} \cdot \delta \sqrt{n}) = \mathcal{O}(\delta n^2)$.
3. Selecting node with the highest potential. This can be done in $\mathcal{O}(n)$.

Adding up the individual terms, we get the complexity of one iteration to be at most $\mathcal{O}(\delta n^2)$. As we aim for the resulting access node set of size $\mathcal{O}(\sqrt{n})$, we would get the total running time of $\mathcal{O}(\delta n^{2.5})$. However, we remind that the algorithm is only a heuristics with no guarantees on the resulting access node set size ²³.

The resulting running time is still quite impractical for bigger timetables. For example, the computation on the dataset *snCF* took more than an hour. This is due to the initial iterations, during which average neighbourhood is still very large (spanning almost the whole graph) and thus in the first two phases we have to do a lot of re-computations. We therefore embrace a simple trick: we do not start with $\mathcal{A} = \emptyset$ but with some access nodes already selected based on high degree. We chose to start with $\frac{\sqrt{n}}{2}$ nodes with the highest degree ($\mathcal{A}_{deg}(\frac{\sqrt{n}}{2})$) - enough to speed-up the computation but not influencing the resulting AN set too much.

The access node sets chosen with the *Locsep* algorithm were much smaller then those selected by the previous approaches and with the increasing number of stations, the r_{locsep}^{avg} value was found to increase only slightly, as could be seen from plots 5.18. The average number of local access nodes for each city was also found to be very small and increasing only very little with increasing n .

To sum up, in *all* of our datasets ²⁴, each scaled to *various* sizes, we were always able to find (r_1, r_2, r_3) access node set \mathcal{A} with the *Locsep* algorithm, such that:

²³The algorithm basically selects access nodes on a greedy basis. However, even that is done only heuristically, using only local scope to reduce the time complexity.

²⁴Except *air01*, which is a special type of timetable.

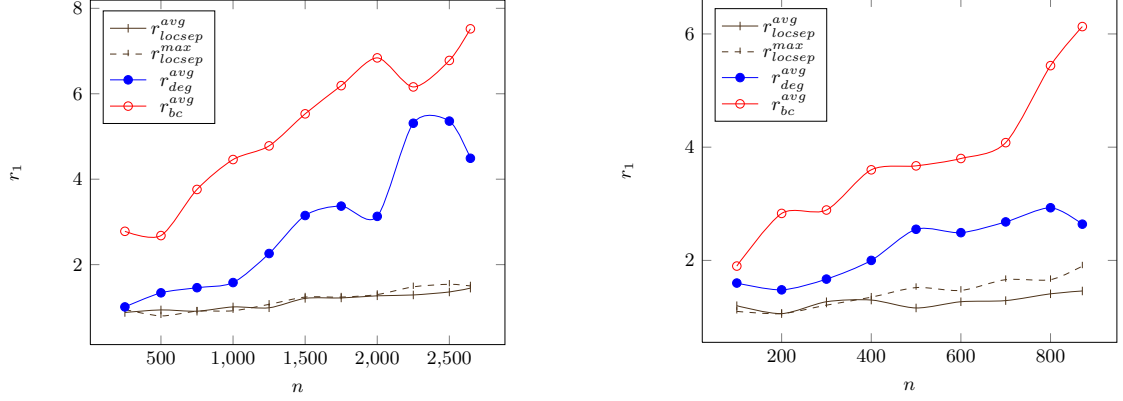


Figure 5.18: The necessary size of \mathcal{A} when selecting ANs based on degree, BC or with *Locsep* algorithm. Datasets *snf* (left) and *cpsk* (right). An ideal situation would be a constant or non-increasing function, to which *Locsep* comes closest.

- $r_1 \leq \text{todo}$
- $r_2 \leq \text{todo}$
- $r_3 \leq \text{todo}$
- \mathcal{A} can be found in $\mathcal{O}(\delta n^{2.5})$

Therefore, when we use *USP-OR-A* together with *Locsep* on our timetables, we achieve parameters as described in table 5.5.

<i>USP-OR-A</i> + <i>Locsep</i>	<i>prep</i>	<i>size</i>	<i>qtime</i>	<i>stretch</i>
Our timetables	$\mathcal{O}(\delta n^{2.5})$	$\mathcal{O}(n^{1.5})$	avg. $\mathcal{O}(\sqrt{n} \log n)$	1

Table 5.5: Parameters for *USP-OR-A* with *Locsep*.

5.4 Performance and comparisons

In this subsection we give the results of the performance of our algorithms on our datasets. We focus on query time and space complexity of the preprocessed oracles. We have already introduced the speed-up as the ratio of average query time for the TD Dijkstra and the average query time for the given algorithm. We will have a similar measure for the size of the preprocessed data, which we compare against the amount of data needed to represent the actual timetable itself.

Definition 5.6. Size-up ($\text{szp}(m)$)

A size-up of an oracle based method m is the ratio $\frac{\text{size}(TD)}{\text{size}(m)}$ where $\text{size}(TD)$ is the size of the memory necessary to store the time-dependent graph.

5.4.1 Performance of *USP-OR*

Query time-wise, *USP-OR* outperforms time-dependent Dijkstra's algorithm almost 80 times (on the subset of *snf-ter* dataset). However, this was at the cost of high space consumption of the method, in some cases requiring almost 400 times more memory than necessary for storage of the time-dependent graph. Therefore, we were not even able to preprocess the method for some of our bigger datasets. Table 5.6 gives a good overview of achieved speed-ups and size-ups.

Name	n	spd	szp
<i>cpru</i> *	700	14.5	396.7
<i>cpza</i> *	700	14.3	265.1
<i>montr</i>	217	8.8	61.1
<i>sncf</i> *	1000	64.8	106.2
<i>sncf-inter</i>	366	27.0	30.3
<i>sncf-ter</i> *	1000	78.3	87.4
<i>zsr</i> (daily)	233	19.3	60.8

Table 5.6: Speed-ups and size-ups of the *USP-OR* algorithm for the whole timetables (for those marked with asterisk we took only a subset of n stations, as we were limited by the space).

In both timetables from *cp.sk* we obtained quite similar results. The query time of *USP-OR* mildly rises with increasing n . This is due to two (not surprising) facts:

- The USP coefficient gets slightly bigger with increasing n
- The OC radius increases too with increasing n

The speed-up was up to 15, however, at the cost of very high space complexity, which made it possible to try out only timetables of the size up to 700. With the *montr* dataset we got similar results, but on a smaller scale.

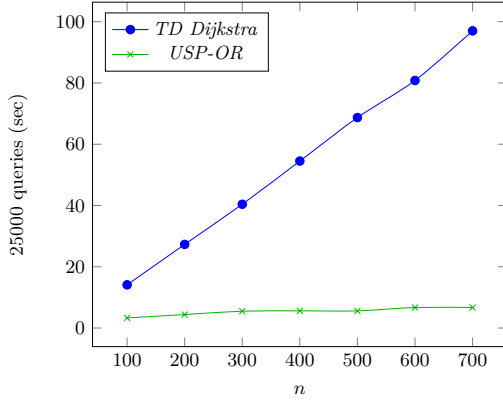


Figure 5.19: **Query time** of *USP-OR* algorithm compared to TD Dijkstra on the *cpru* dataset. **Changing n .**

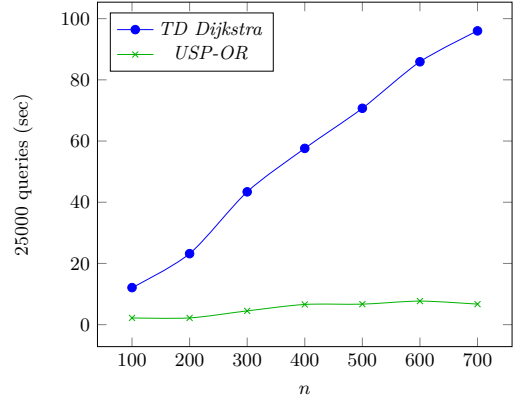


Figure 5.20: **Query time** of *USP-OR* algorithm compared to TD Dijkstra on the *cpza* dataset. **Changing n .**

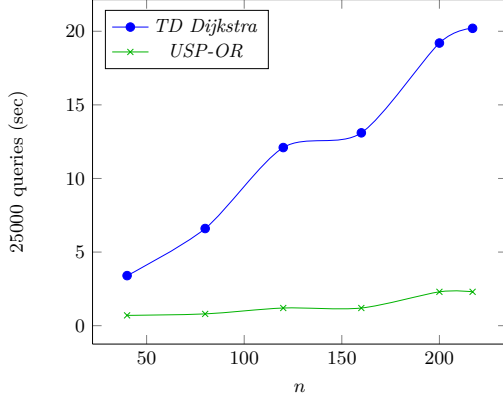


Figure 5.21: **Query time** of *USP-OR* algorithm compared to TD Dijkstra on the *montr* dataset. **Changing n .**

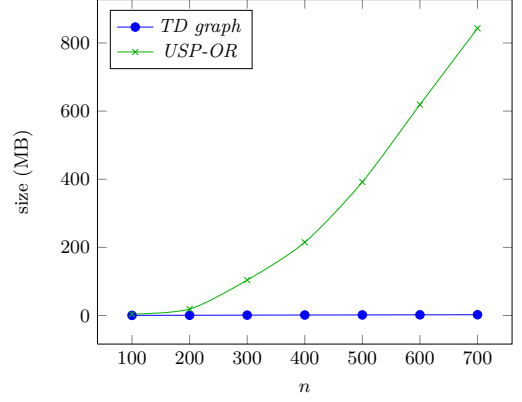


Figure 5.22: **Size** (in MB) of the oracle for *USP-OR* vs. size of TD graph on *cpza* dataset. **Changing n .**

In the datasets from SNCF, an interesting thing was that the the query-time actually decreased with increased size. This was due to the average USP coefficient getting smaller in bigger datasets while OC radius not increasing too much. We measured here the speed-ups of up to 80 for *snCF-ter*, with smaller size-ups then in case of *cp.sk* timetables.

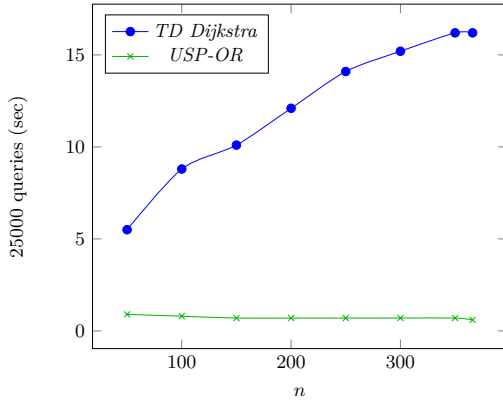


Figure 5.23: **Query time** of *USP-OR* algorithm compared to TD Dijkstra on the *snCF-inter* dataset. **Changing n .**

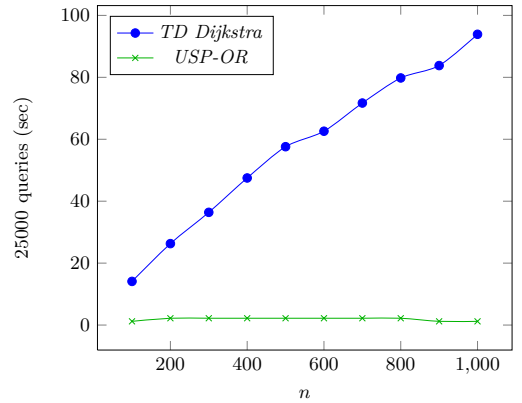


Figure 5.24: **Query time** of *USP-OR* algorithm compared to TD Dijkstra on the *snCF-ter* dataset. **Changing n .**

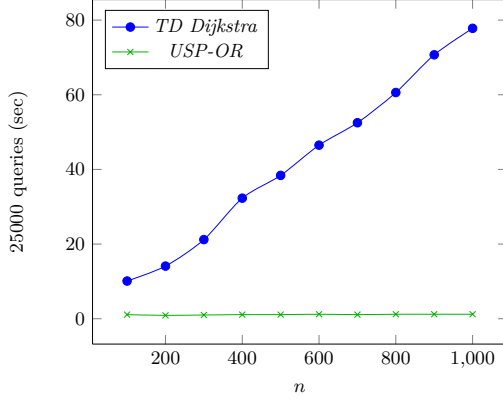


Figure 5.25: **Query time** of *USP-OR* algorithm compared to TD Dijkstra on the *snCF* dataset. **Changing n .**

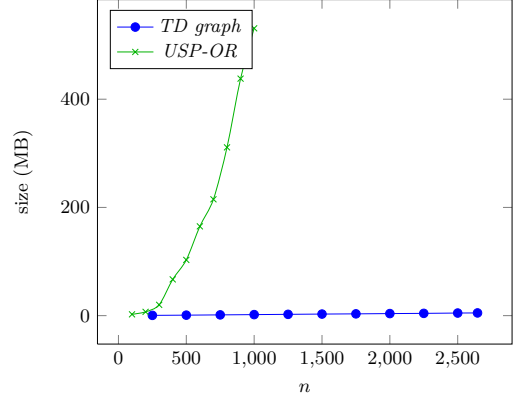


Figure 5.26: **Size** (in MB) of the oracle for *USP-OR* vs. size of TD graph on *snCF* dataset. **Changing n .**

On the *zsr* dataset we measured how increased time range influences the query time. You may see that for both algorithms the query time almost stops increasing at some point - this is because (informally) adding time range no longer brings along new optimal connections (or underlying shortest paths in case of *USP-OR*).

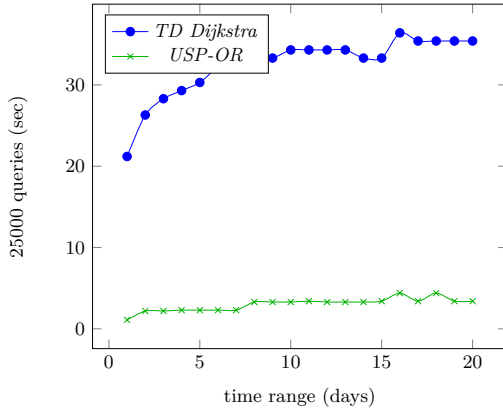


Figure 5.27: **Query time** of *USP-OR* algorithm compared to TD Dijkstra on the *zsr* dataset. **Changing r .**

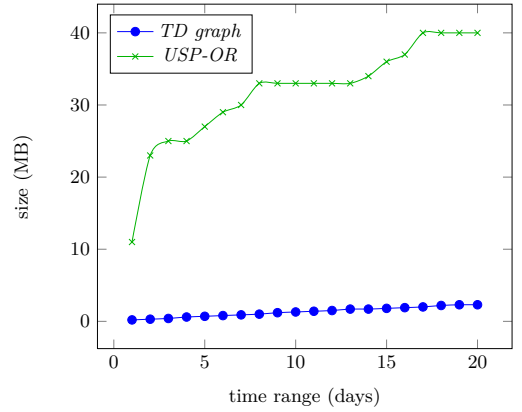


Figure 5.28: **Size** (in MB) of the oracle for *USP-OR* vs. size of TD graph on *zsr* dataset. **Changing r .**

5.4.2 *USP-OR-A* with *Locsep*

With *USP-OR-A*, the speed-ups were no longer so high as in case of *USP-OR-A*, but neither were the size-ups, so we could fully try out all of our datasets. The maximum speed-up was achieved in *snCF-ter*, where the *USP-OR-A* outperformed *TD Dijkstra* almost 7 times. In our datasets, the preprocessed oracle of *USP-OR-A* did not need more than 3 times the size of the TD graph. Table 5.7 provides an overview of achieved speed-ups and size-ups.

Name	n	spd	szp
<i>cpru</i>	871	1.8	2.97
<i>cpza</i>	1108	1.9	2.52
<i>montr</i>	217	1.6	1.14
<i>sncf</i>	2646	6.3	3.0
<i>sncf-inter</i>	366	3.9	1.59
<i>sncf-ter</i>	2637	6.9	2.43
<i>zsr</i> (daily)	233	2.5	1.99

Table 5.7: Speed-ups and size-ups of the *USP-OR-A* with *Locsep* for the whole timetables (for those marked with asterisk we took only a subset of n stations, as we were limited by the space).

The bus timetables proved to be a bigger challenge for *USP-OR-A*, achieving milder speed-ups and requiring more memory then railways timetables. However, bigger timetables would be necessary to obtain more relevant results.

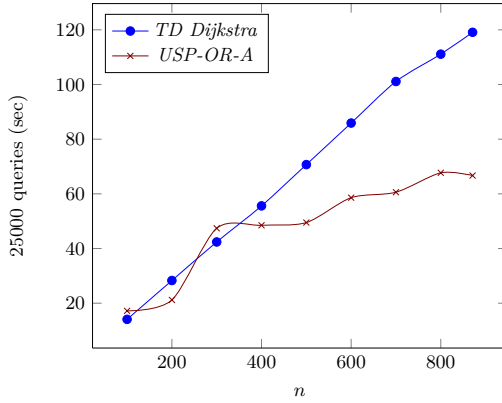


Figure 5.29: **Query time** of *USP-OR-A* with *Locsep* compared to TD Dijkstra on the *cpru* dataset. **Changing n .**

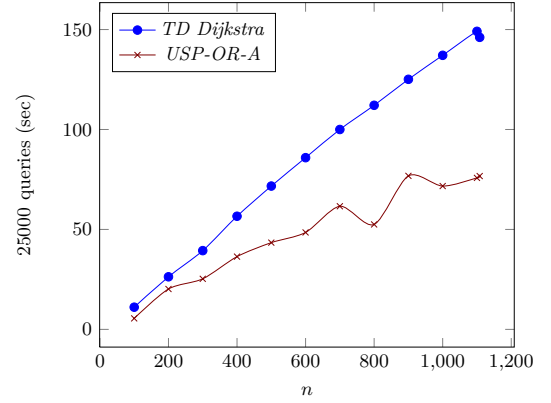


Figure 5.30: **Query time** of *USP-OR-A* with *Locsep* compared to TD Dijkstra on the *cpza* dataset. **Changing n .**

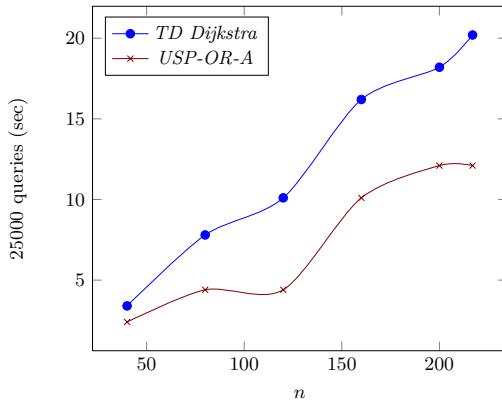


Figure 5.31: **Query time** of *USP-OR-A* with *Locsep* compared to TD Dijkstra on the *montr* dataset. **Changing n .**

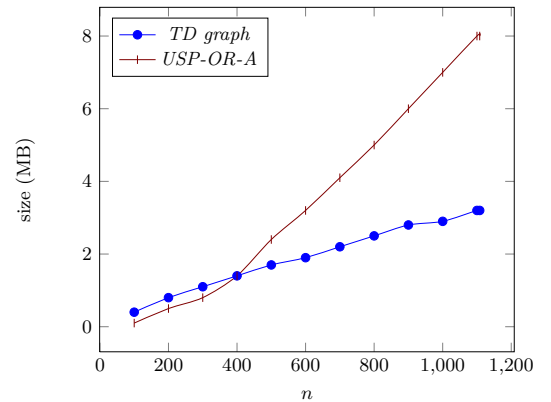


Figure 5.32: **Size** (in MB) of the oracle for *USP-OR-A* with *Locsep* vs. size of TD graph on *cpza* dataset. **Changing n .**

In our biggest datasets, we achieved the best speed-ups while the size-up still stayed relatively small (though here we better see its tendency to increase as $n^{1.5}$). It would be interesting to try out even bigger datasets as the speed-up was gradually increasing with increasing n .

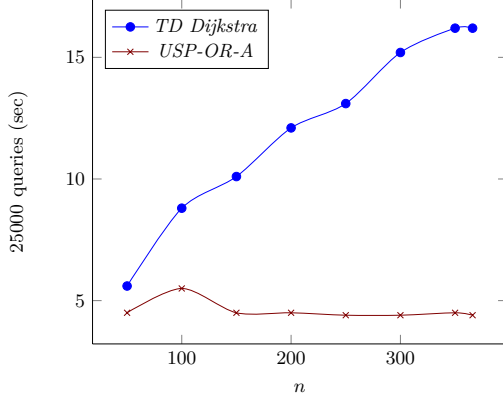


Figure 5.33: **Query time** of *USP-OR-A* with *Locsep* compared to TD Dijkstra on the *snCF-inter* dataset. **Changing n .**

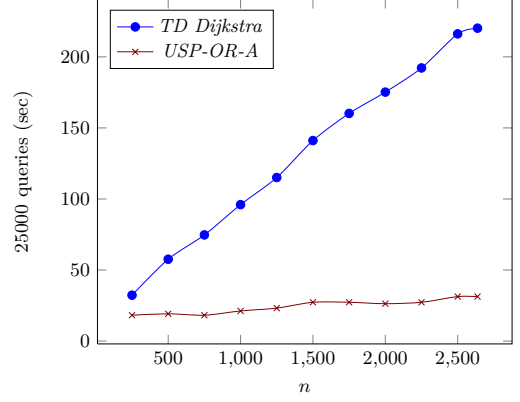


Figure 5.34: **Query time** of *USP-OR-A* with *Locsep* compared to TD Dijkstra on the *snCF-ter* dataset. **Changing n .**

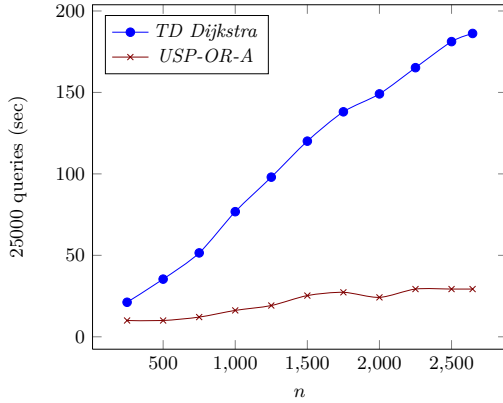


Figure 5.35: **Query time** of *USP-OR-A* with *Locsep* compared to TD Dijkstra on the *snCF* dataset. **Changing n .**

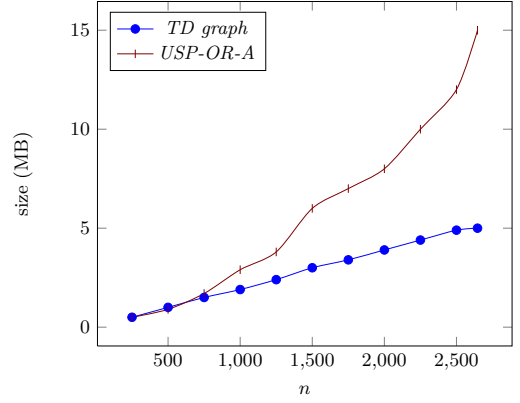


Figure 5.36: **Size (in MB)** of the oracle for *USP-OR-A* with *Locsep* vs. size of TD graph on *snCF* dataset. **Changing n .**

Finally, on the *zsr* timetable, we see two things:

- The space-complexity of *USP-OR-A* is left pretty much unaffected with increased time range.
- The speed-up decreases (since with increased time range, there are generally more USPs between pairs of cities which we have to try out during the query)

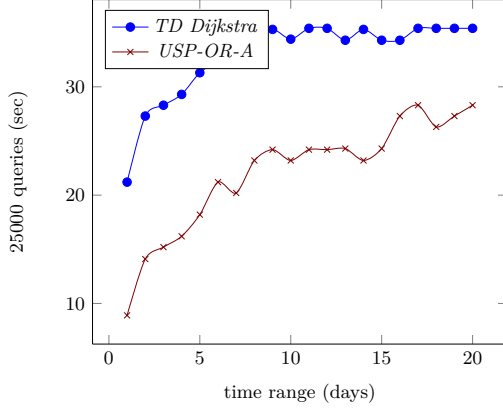


Figure 5.37: **Query time** of *USP-OR-A* with *Locsep* compared to TD Dijkstra on the *zsr* dataset. **Changing r .**

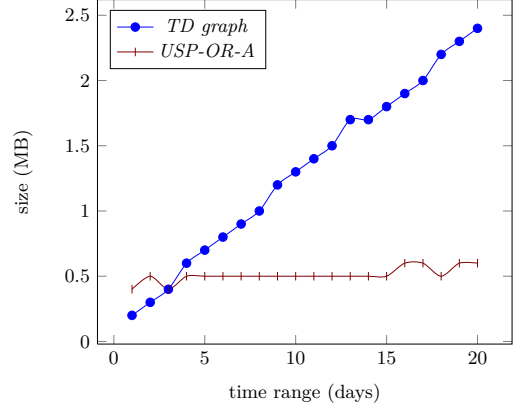


Figure 5.38: **Size** (in MB) of the oracle for *USP-OR-A* with *Locsep* vs. size of TD graph on *zsr* dataset. **Changing r .**

5.4.3 *USP-OR-A* with *Locsep Max*

We also tried out *USP-OR-A* with *Locsep Max* to see if the difference in the stopping criterion of *Locsep* would influence the query times. It did help, but the difference in the performance is minimal, therefore we list only the table summarizing the speed-ups and size-ups (5.8) and the details for datasets *cpza* and *sncf*.

Name	n	spd	szp
<i>cpru</i>	871	2.1	4.5
<i>cpza</i>	1108	2.1	3.1
<i>montr</i>	217	2.1	1.9
<i>sncf</i>	2646	6.6	3.0
<i>sncf-inter</i>	366	4.3	1.6
<i>sncf-ter</i>	2637	7.1	2.43
<i>zsr</i> (daily)	233	2.3	1.94

Table 5.8: Speed-ups and size-ups of the *USP-OR-A* with *Locsep Max* for the whole timetables (for those marked with asterisk we took only a subset of n stations, as we were limited by the space).

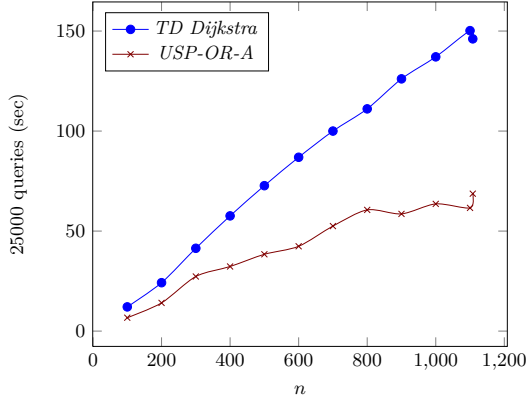


Figure 5.39: **Query time** of *USP-OR-A* with *Locsep Max* compared to TD Dijkstra on the *cpza* dataset. **Changing n .**

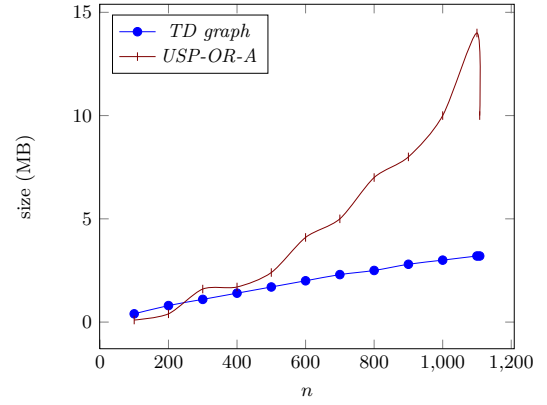


Figure 5.40: **Size** (in MB) of the oracle for *USP-OR-A* with *Locsep Max* vs. size of TD graph on *cpza* dataset. **Changing n .**

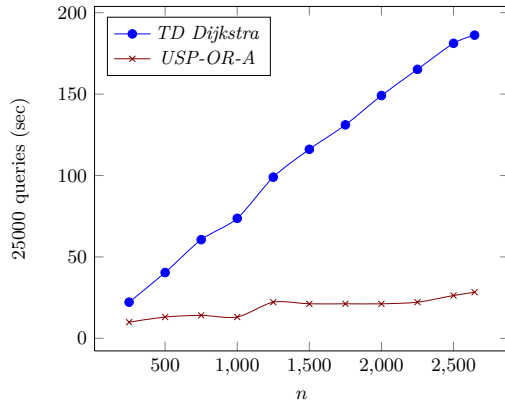


Figure 5.41: **Query time** of *USP-OR-A* with *Locsep Max* compared to TD Dijkstra on the *sncf* dataset. **Changing n .**

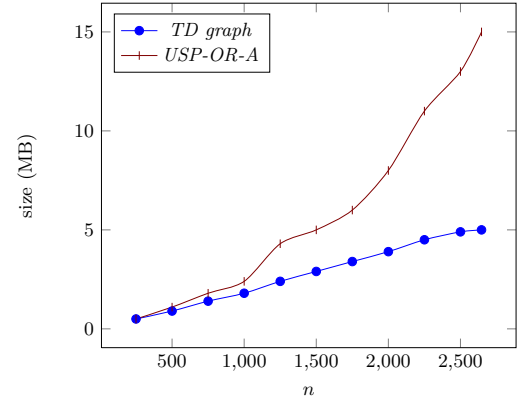


Figure 5.42: **Size** (in MB) of the oracle for *USP-OR-A* with *Locsep Max* vs. size of TD graph on *sncf* dataset. **Changing n .**

6 Neural network approach

7 Application TTBlazer

8 Conclusion

Appendices

A File formats

Timetable is simply a set of elementary connections, thus the format is:

- number of el. connections
- the list of all el. connections (one per line, format “*FROM TO DEP-DAY DEP-TIME ARR-DAY ARR-TIME*”)

```
1 7 //number of elementary connections
2 A B 0 10:00 0 10:45 //el. connection
3 A B 0 11:00 0 11:45
4 A B 0 12:00 0 12:45
5 A C 0 09:30 0 10:00
6 A C 0 10:15 0 10:45
7 C D 0 11:00 0 11:30
8 C D 0 13:00 0 13:30
```

Listing 1: TT file format.

Underlying graph is basically an oriented graph, with some optional parameters. The format is the following:

- number of cities
- number of arcs
- the list of all cities (one per line)
 - optional coordinates (otherwise null)
- the list of all arcs (one per line, format “*FROM TO*”)
 - optional length (otherwise null)
 - optional list of lines operating on that arc (otherwise null)

```
1 4 //number of cities
2 5 //number of arcs
3 A 45 32 //name of the city, optional coordinates
4 B null
5 C 56 34
6 D null
7 A B 57 Northern //arc, optional length and list of lines
8 A C null Picadilly Victoria
9 C B 45 Circle Jubilee Picadilly
10 C D 32 null
11 D A null null
```

Listing 2: UG file format.

Time-expanded graph is simply an oriented weighted graph, with nodes being the events and arcs being the elementary connections or waiting edges:

- number of nodes (i.e. events)
- number of arcs (el. connections + waiting)
- the list of all events (in the format “*CITY DAY TIME*”)
- the list of all arcs (in the format “*FROM-EVENT TO-EVENT*”)

```

1 5 //number of events
2 15 //number of arcs
3 A 0 13:30 //event
4 A 0 14:00
5 B 0 13:45
6 B 0 15:00
7 C 0 14:15
8 A 0 13:30 A 0 14:00 //waiting arc
9 A 0 13:30 B 0 13:45 //el. connection arc
10 A 0 14:00 B 0 15:00
11 A 0 13:30 B 0 15:00
12 C 0 14:15 B 0 15:00
13 ...

```

Listing 3: TE file format.

Time-dependent graph is an oriented graph with a function on the arc specifying the arc's traversal time at any moment. In timetable networks this function is piece-wise linear and it is fully represented by the list of its interpolation points. Thus the TD file format:

- number of cities
- number of arcs
- the list of all cities (one per line)
 - optional coordinates (otherwise null)
- the list of all arcs (one per line). Arc has the format “*FROM TO INT-POINTS*” where *INT-POINTS* is a list of interpolation points²⁵, see the listing 4 for an example.

```

1 4 //number of stations
2 5 //number of arcs
3 A 0 0 //name of the city, optional coordinates
4 B 4 4
5 C null
6 D 12 0
7 A B (0 13:30 45) (0 14:00 40) //arc and the list of interpolation
    points
8 A C (1 14:15 10)
9 C B (0 15:00 20)
10 C D (2 10:00 70)
11 D A (1 17:20 35) (1 18:00 40) (1 18:50 35)
12 ...

```

Listing 4: TD file format.

²⁵An interpolation point is described by a triple “*DAY TIME MINUTES*”, where *MINUTES* are the traversal time