

Multi-Modal Route Planning in Road and Transit Networks

Master's Thesis

Daniel Tischner

University of Freiburg, Germany,
`daniel.tischner.cs@gmail.com`

August 12, 2018

Supervisor: Prof. Dr. Hannah Bast
Advisor: Patrick Brosi

Contents

1	Introduction	6
2	Preliminaries	7
2.1	Graph	7
2.2	Tree	9
2.3	Metric	9
3	Models	12
3.1	Road graph	12
3.2	Transit graph	14
3.3	Link graph	17
3.4	Timetable	19
4	Nearest neighbor problem	21
4.1	Cover tree	23
5	Shortest path problem	28
5.1	Time-independent	29
5.1.1	Dijkstra	30
5.1.2	A* and ALT	32
5.2	Time-dependent	34
5.2.1	Connection scan	34
5.3	Multi-modal	38
5.3.1	Modified Dijkstra	38
5.3.2	Access nodes	38
5.4	Other algorithms	38
6	Evaluation	38
6.1	Input data	38
6.2	Experiments	38
6.2.1	Nearest neighbor computation	38
6.2.2	Uni-modal routing	39
6.2.3	Multi-modal routing	39
6.3	Summary	39
7	Conclusion	39
	References	39

Declaration

I hereby declare, that I am the sole author and composer of my Thesis and that no other sources or learning aids, other than those listed, have been used. Furthermore, I declare that I have acknowledged the work of others by providing detailed references of said work. I hereby also declare, that my Thesis has not been prepared for another examination or assignment, either wholly or excerpts thereof.

Place, Date

Signature

Zusammenfassung

Wir präsentieren Algorithmen für multi-modale Routenplanung in Straßennetzen und Netzwerken des öffentlichen Personennahverkehrs (ÖPNV), so wie in kombinierten Netzwerken.

Dazu stellen wir das Nächste-Nachbar und das Kürzester-Pfad Problem vor und schlagen Lösungen basierend auf COVER TREES, ALT und CSA vor.

Des Weiteren erläutern wir die Theorie hinter den Algorithmen, geben eine kurze Übersicht über andere Techniken, zeigen Versuchsergebnisse auf und vergleichen die Techniken untereinander.

Abstract

We present algorithms for multi-modal route planning in road and public transit networks, as well as in combined networks.

Therefore, we explore the nearest neighbor and shortest path problem and propose solutions based on **COVER TREES**, **ALT** and **CSA**.

Further, we illustrate the theory behind the algorithms, give a short overview of other techniques, present experimental results and compare the techniques with each other.

TODO: disable todos and macro highlights.

1 Introduction

Route planning refers to the problem of finding an *optimal* route between given locations in a network. With the ongoing expansion of road and public transit networks all over the world route planner gain more and more importance. This led to a rapid increase in research [4, 10, 18] of relevant topics and development of route planner software [14, 13, 24].

However, a common problem of most such services is that they are limited to one transportation mode only. That is a route can only be taken by a car or train but not by both at the same time. This is known as uni-modal routing. In contrast to that multi-modal routing allows the alternation of transportation modes. For example a route that first uses a car to drive to a train station, then a train which travels to a another train station and finally using a bicycle from there to reach the destination.

The difficulty with multi-modal routing lies in most algorithms being fitted to networks with specific properties. Unfortunately, road networks differ a lot from public transit networks. As such, a route planning algorithm fitted to a certain type of network will likely yield undesired results, have an impractical running time or not even be able to be used at all on different networks. We will explore this later in **Section 6**.

In this thesis we explore a technique with which we can combine an algorithm fitted for road networks with an algorithm for public transit networks. Effectively obtaining a generic algorithm that is able to compute routes on combined networks. The basic idea is simple, given a source and destination, both in the road network, we select *access nodes* for both. These are nodes where we will switch from the road into the public transit network. A route can then be computed by using the road algorithm for the source to its access nodes, the transit algorithm for the access nodes of the source to the access nodes of the destination and finally the road algorithm again for the destinations access nodes to the destination. Note that this technique might not yield the shortest possible path anymore. Also, it does not allow an arbitrary alternation of transportation modes. However, we accept those limitations since the resulting algorithm is very generic and able to compute routes faster than without limitations. We will cover this technique in detail in **Section 5.3.2**.

Our final technique uses a modified version of **ALT** [15] as road algorithm and **CSA** [11] for the transportation network. The algorithms are presented in **Section 5.1.2** and **Section 5.2.1** respectively. We also develop a multi-modal variant of **DIJKSTRA** [8] which is able to compute the shortest route in a combined network with the possibility of changing transportation modes arbitrarily. It is presented in **Section 5.3.1** and acts as baseline to our final technique based on access nodes.

We compute access nodes by solving the **NEAREST NEIGHBOR PROBLEM**. For a given node in the road network its access nodes are then all nodes in the transit network which are in the *vicinity* of the road node. We explore a solution to this problem in **Section 4**.

Section 3 starts by defining types of networks. We represent road networks by graphs only. For transit networks we provide a graph representation too. Both graphs can then be combined into a linked graph. The advantage of graph based models is that they are well studied and therefore we are able to use our multi-modal variant of **DIJKSTRA** to compute routes on them. However, we also propose a non-graph based representation for transit networks, a timetable. The timetable is used by **CSA**, an efficient algorithm for route planning on public transit networks. With that, our road and transit networks get incompatible and can not easily be combined. Therefore, we use the previously mentioned generic approach based on access nodes for this type of network.

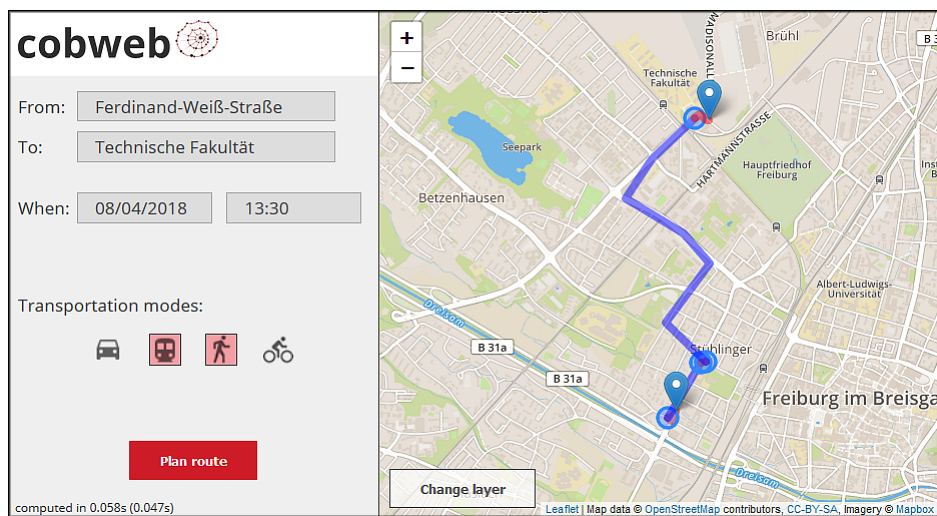


Fig. 1: Screenshot of **COBWEBs** [22] frontend, an open-source multi-modal route planner. It shows a multi-modal route starting from a given source, using the modes *foot-tram-foot-tram-foot* in that sequence to reach the destination.

Further, we implemented the presented algorithms in the **COBWEB** [22] project, which is an open-source multi-modal route planner (see **Fig. 1** for an image of its frontend). In **Section 6** we show our experimental results and compare the techniques with each other.

2 Preliminaries

Before we define our specific data models and problems we will introduce and formalize commonly reoccurring terms.

2.1 Graph

Definition 1. A graph G is a tuple (V, E) with a set of nodes V and a set of edges $E \subseteq V \times \mathbb{R}_{\geq 0} \times V$. An edge $e \in E$ is an ordered tuple (u, w, v) with source node $u \in V$,

a non-negative weight $w \in \mathbb{R}_{\geq 0}$ and a destination node $v \in V$.

Note that **Definition 1** actually defines a *directed* graph, as opposed to an *undirected* graph where an edge like (u, w, v) would be considered equal to the edge of opposite direction (v, w, u) (compare to [12]). However, for transportation networks an undirected graph often is not applicable, for example due to one way streets or time dependent connections like trains which depart at different times for different directions.

In the context of route planning we refer to the weight w of an edge (u, w, v) as *cost*. It can be used to encode the length of the represented connection. Or to represent the time it takes to travel the distance in a given transportation mode.

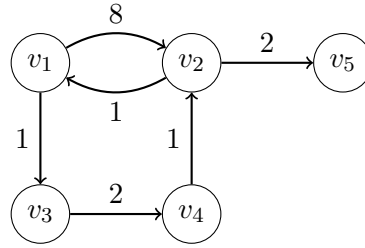


Fig. 2: Illustration of an example graph with five nodes and six edges.

As an example consider the graph $G = (V, E)$ with

$$V = \{v_1, v_2, v_3, v_4, v_5\} \text{ and}$$

$$E = \{(v_1, 8, v_2), (v_1, 1, v_3), (v_2, 1, v_1), (v_2, 2, v_5), (v_3, 2, v_4), (v_4, 1, v_2)\}.$$

which is illustrated by **Fig. 2**.

Definition 2. Given a graph $G = (V, E)$ the function $\text{src} : E \rightarrow V, ((u, w, v)) \mapsto u$ gets the source of an edge. Analogously $\text{dest} : E \rightarrow V, ((u, w, v)) \mapsto v$ retrieves the destination.

Definition 3. A path in a graph $G = (V, E)$ is a sequence $p = e_1 e_2 e_3 \dots$ of edges $e_i \in E$ such that

$$\forall i : \text{dest}(e_i) = \text{src}(e_{i+1}).$$

The length of a path is the amount of edges it contains, i.e. the length of the sequence. The weight or cost is the sum of its edges weights.

An example for a path in the graph G would be

$$p = (v_1, 8, v_2)(v_2, 1, v_1)(v_1, 1, v_3).$$

Its length is 3 and it has a weight of 10.

2.2 Tree

Definition 4. A tree is a graph $T = (V, E)$ with the following properties:

1. There is exactly one node $r \in V$ with no ingoing edges, called the root, i.e.

$$\exists! r \in V \nexists e \in E : \text{dest}(e) = r.$$

2. All other nodes v have exactly one ingoing edge. The source p of this edge is called parent of v and v is called child of p :

$$\forall v \in V : v \neq r \Rightarrow \exists! e \in E : \text{dest}(e) = v.$$

Definition 5. The subtree of a tree $T = (V, E)$ rooted at a node $r' \in V$ is a tree $T' = (V', E')$. $V' \subseteq V$ is the set of nodes that can be reached from r' . That is, all nodes that are part of possible paths starting at r' . Likewise $E' \subseteq E$ is the set of edges restricted to the vertices in V' . The root of T' is r' .

Definition 6. The depth of a node v in a tree $T = (V, E)$, denoted by $\text{depth}(v)$, is defined as the amount of edges between the v and the root r . It is the length of the unique path p starting at r and ending at v .

The height of a tree is its greatest depth, i.e.

$$\max_{v \in V} \text{depth}(v).$$

And

$$\text{children}(v) = \{c \in T \mid c \text{ child of } v\}.$$

Trees are hierarchical data structures. Every node, except the root, has one parent. A node itself can have multiple children. Note that it is not possible to form a loop in a tree, i.e. a path that visits a node more than once. A node without children is called a leaf. TODO: Maybe add some illustrations here...

2.3 Metric

Definition 7. A function $d : M \times M \rightarrow \mathbb{R}$ on a set M is called a metric iff for all $x, y, z \in M$

$d(x, y) \geq 0,$	<i>non-negativity</i>
$d(x, y) = 0 \Leftrightarrow x = y,$	<i>identity of indiscernibles</i>
$d(x, y) = d(y, x)$ and	<i>symmetry</i>
$d(x, z) \leq d(x, y) + d(y, z)$	<i>triangle inequality</i>

holds.

Definition 8. A metric space is a pair (M, d) where M is a set and $d : M \times M \rightarrow \mathbb{R}$ a metric on M .

Definition 9. Given a metric d on a set M , the distance of a point $p \in M$ to a subset $Q \subseteq M$ is defined as the distance from p to its nearest point in Q :

$$d(p, Q) = \min_{q \in Q} d(p, q)$$

A metric is used to measure the distance between given locations. **Section 4** and **Section 5**, in particular **Section 5.1.2**, will make heavy use of this term.

There we measure the distance between geographical locations given as pair of *latitude* and *longitude* coordinates. Latitude and longitude, often denoted by ϕ and λ , are real numbers in the ranges $(-90, 90)$ and $[-180, 180)$ respectively, measured in degrees. However, for convenience we represent them in radians. Both representations are equivalent to each other and can easily be converted using the ratio $360^\circ = 2\pi$ rad.

A commonly used measure is the *as-the-crow-flies* metric, which is equivalent to the euclidean distance in the euclidean space. **Definition 10** defines an approximation of this distance on locations given by latitude and longitude coordinates. The approximation is commonly known as equirectangular projection of the earth [19]. Note that there are more accurate methods for computing the great-circle distance for geographical locations, like the haversine formula [20]. However, they come with a significant computational overhead.

Definition 10. Given a set of coordinates $M = \{(\phi, \lambda) | \phi \in (-\frac{\pi}{2}, \frac{\pi}{2}), \lambda \in [-\pi, \pi)\}$ we define $\text{asTheCrowFlies} : M \times M \rightarrow \mathbb{R}$ such that

$$((\phi_1, \lambda_1), (\phi_2, \lambda_2)) \mapsto \sqrt{\left((\lambda_2 - \lambda_1) \cdot \cos\left(\frac{\phi_1 + \phi_2}{2}\right)\right)^2 + (\phi_2 - \phi_1)^2 \cdot 6371000}.$$

As a next step we prove that asTheCrowFlies is indeed a metric on the set of coordinates.

Proposition 1. The function asTheCrowFlies is a metric on its domain M .

Proof. We need to prove that all four axioms hold. Let us first set

$$\begin{aligned} x &= (\lambda_2 - \lambda_1) \cdot \cos\left(\frac{\phi_1 + \phi_2}{2}\right) \\ y &= \phi_2 - \phi_1 \end{aligned}$$

then the function simplifies to

$$\sqrt{x^2 + y^2 \cdot 6371000}.$$

Obviously this can never resolve to a negative number since

$$\underbrace{\underbrace{\underbrace{x^2}_{\geq 0} + \underbrace{y^2}_{\geq 0}}_{\geq 0}}_{\geq 0} \cdot 6371000.$$

For the second axiom we assume that $\text{asTheCrowFlies}((\phi_1, \lambda_1), (\phi_2, \lambda_2)) = 0$ for an arbitrary pair of coordinates and follow

$$\begin{aligned} & \sqrt{\left((\lambda_2 - \lambda_1) \cdot \cos\left(\frac{\phi_1 + \phi_2}{2}\right)\right)^2 + (\phi_2 - \phi_1)^2} \cdot 6371000 = 0 \\ \Leftrightarrow & \sqrt{\left((\lambda_2 - \lambda_1) \cdot \cos\left(\frac{\phi_1 + \phi_2}{2}\right)\right)^2 + (\phi_2 - \phi_1)^2} = 0 \\ \Leftrightarrow & \left((\lambda_2 - \lambda_1) \cdot \cos\left(\frac{\phi_1 + \phi_2}{2}\right)\right)^2 + (\phi_2 - \phi_1)^2 = 0 \end{aligned}$$

At this point either both summands are 0 or one is the negative of the other. However, both summands must be positive due to the quadration. Because of that we follow

$$\begin{aligned} & (\phi_2 - \phi_1)^2 = 0 \\ \Leftrightarrow & \phi_2 = \phi_1 \end{aligned}$$

and with that

$$\begin{aligned} & \left((\lambda_2 - \lambda_1) \cdot \cos\left(\frac{\phi_1 + \phi_2}{2}\right)\right)^2 = 0 \\ \Leftrightarrow & (\lambda_2 - \lambda_1) \cdot \cos\left(\frac{2\phi_1}{2}\right) = 0 \\ \Leftrightarrow & (\lambda_2 - \lambda_1) \cdot \cos(\phi_1) = 0. \end{aligned}$$

Since $\phi_1 \in (-\frac{\pi}{2}, \frac{\pi}{2})$ it follows that $\cos(\phi_1) \neq 0$. As such

$$\begin{aligned} & \lambda_2 - \lambda_1 = 0 \\ \Leftrightarrow & \lambda_2 = \lambda_1 \end{aligned}$$

and by that $(\phi_1, \lambda_1) = (\phi_2, \lambda_2)$, so the second axiom holds.

Symmetry follows quickly since

$$\begin{aligned} & \phi_1 + \phi_2 = \phi_2 + \phi_1 \\ & (\phi_2 - \phi_1)^2 = (\phi_1 - \phi_2)^2 \\ & \left((\lambda_2 - \lambda_1) \cdot \cos\left(\frac{\phi_1 + \phi_2}{2}\right)\right)^2 = (\lambda_2 - \lambda_1)^2 \cdot \cos^2\left(\frac{\phi_1 + \phi_2}{2}\right) \\ & (\lambda_2 - \lambda_1)^2 = (\lambda_1 - \lambda_2)^2. \end{aligned}$$

The triangle inequality is a bit trickier, we choose three arbitrary coordinates $c_i = (\phi_i, \lambda_i)$ for $i = 1, 2, 3$ and start on the squared left side:

$$\begin{aligned}
\text{asTheCrowFlies}^2(c_1, c_3) &= \left(\left((\lambda_3 - \lambda_1) \cdot \cos\left(\frac{\phi_1 + \phi_3}{2}\right) \right)^2 + (\phi_3 - \phi_1)^2 \right) \cdot 6371000^2 \\
&= \left(\left((\lambda_3 - \lambda_2 + \lambda_2 - \lambda_1) \cdot \cos\left(\frac{\phi_1 + \phi_3}{2}\right) \right)^2 + (\phi_3 - \phi_2 + \phi_2 - \phi_1)^2 \right) \cdot 6371000^2 \\
&= \left(\left((\lambda_3 - \lambda_2)^2 + (\lambda_2 - \lambda_1)^2 + 2 \cdot ((\lambda_3 - \lambda_2) \cdot (\lambda_2 - \lambda_1)) \right) \cdot \cos^2\left(\frac{\phi_1 + \phi_3}{2}\right) \right. \\
&\quad \left. + (\phi_3 - \phi_2)^2 + (\phi_2 - \phi_1)^2 + 2 \cdot ((\phi_3 - \phi_2) \cdot (\phi_2 - \phi_1)) \right) \cdot 6371000^2 \\
&= \dots \\
&\leq \left(\left(\left((\lambda_2 - \lambda_1) \cdot \cos\left(\frac{\phi_1 + \phi_2}{2}\right) \right)^2 + (\phi_2 - \phi_1)^2 \right) \cdot 6371000^2 \right) \\
&\quad + \left(\left(\left((\lambda_3 - \lambda_2) \cdot \cos\left(\frac{\phi_2 + \phi_3}{2}\right) \right)^2 + (\phi_3 - \phi_2)^2 \right) \cdot 6371000^2 \right) \\
&\quad + 2 \cdot \left(\left(\left((\lambda_2 - \lambda_1) \cdot \cos\left(\frac{\phi_1 + \phi_2}{2}\right) \right)^2 + (\phi_2 - \phi_1)^2 \right) \cdot 6371000^2 \right) \\
&\quad \cdot \left(\left(\left((\lambda_3 - \lambda_2) \cdot \cos\left(\frac{\phi_2 + \phi_3}{2}\right) \right)^2 + (\phi_3 - \phi_2)^2 \right) \cdot 6371000^2 \right) \\
&= (\text{asTheCrowFlies}(c_1, c_2) + \text{asTheCrowFlies}(c_2, c_3))^2
\end{aligned}$$

TODO: continue (squared ineq holds without, cauchy schwarz ineq) or remove completely...

All four axioms hold, asTheCrowFlies is a metric on the set M . \square

3 Models

This section defines the models we use for the different network types. We define a graph based representation for road and transit networks. Then both graphs are combined into a linked graph, making it possible to have one graph for the whole network. Afterwards an alternative representation for transit networks is shown.

3.1 Road graph

A road network typically is time independent. It consists of geographical locations and roads connecting them with each other. We assume that a road can be taken at any time, with no time dependent constraints (see **Section 2** of [10]).

Modeling the network as graph is straightforward, **Definition 11** goes into detail.

Definition 11. A road graph is a graph $G = (V, E)$ with a set of geographic coordinates

$$V = \{(\phi, \lambda) | \phi \in \left(-\frac{\pi}{2}, \frac{\pi}{2}\right), \lambda \in [-\pi, \pi)\},$$

for example road junctions. There is an edge $(u, w, v) \in E$ iff there is a road connecting the location u with the location v , which can be taken in that direction. The weight w of the edge is the average time needed to take the road from u to v using a car, measured in seconds.

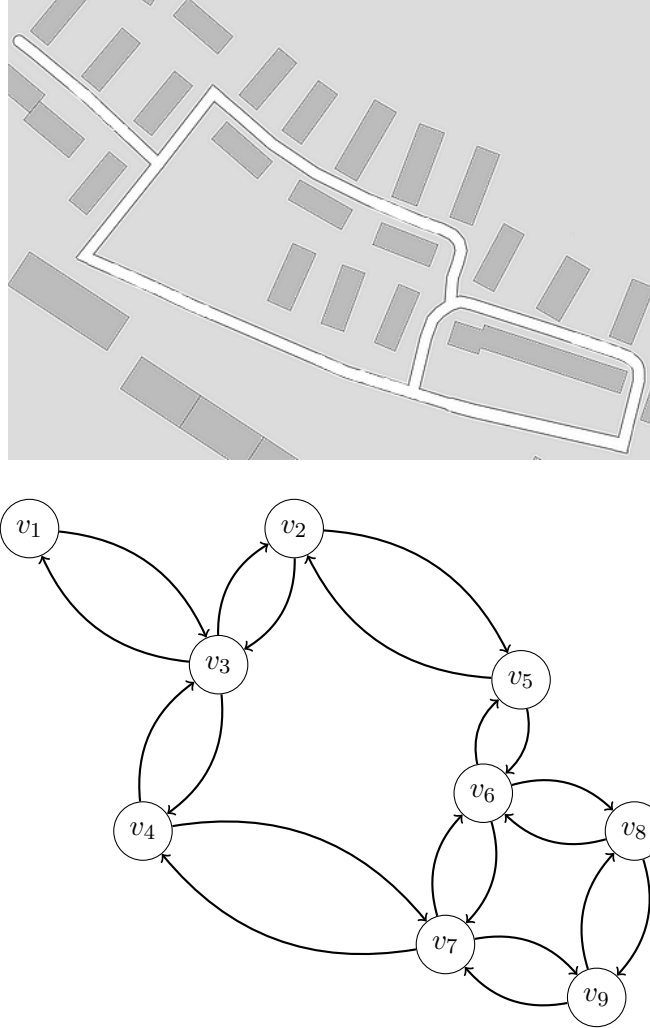


Fig. 3: Example of a road network with its corresponding road graph. White connections indicate roads, dark gray rectangles represent houses or other static objects. Geographical coordinates for each node as well as edge weights are omitted in the graph illustration.

Fig. 3 shows a constructed example road network with the corresponding road graph. Note that two way streets result in two edges, one edge for every direction the road can be taken.

Since edge weights are represented as average time it needs to take the road, it is possible to encode different road types. For example the average speed on a motorway is much higher than on a residential street. As such, the weight of an edge representing a motorway is much smaller than the weight of an edge representing a residential street.

While the example has exactly one node per road junction this must not always be the case. Typical real world data often consists of multiple nodes per road segment. However, **Definition 11** is still valid for such data as long as there are edges between the nodes if and only if there is a road connecting the locations.

3.2 Transit graph

Transit networks can be modeled similar to road graphs. The key difference is that transit networks are time dependent while road networks typically are not. For example an edge connecting *Freiburg main station* with *Karlsruhe main station* can not be taken at any time since trains and other transit vehicles only depart at certain times. The schedule might even change at different days.

The difficulty lies in modeling time dependence in a static graph. There are two common approaches to that problem (see [10, 18, 4]).

The first approach is called *time-dependent*. There edge weights are not static numbers but functions that take a date with time and compute the cost it needs to take the edge when starting at the given time. This includes waiting time. As an example assume an edge (u, c, v) with the cost function c . The edge represents a train connection and the travel time are 10 minutes. However, the train departs at 10:15 *am* but the starting time is 10:00 *am*. The cost function thus computes a waiting time of 15 minutes plus the travel time of 10 minutes. Resulting in an edge weight of 25 minutes.

The main problem with this model is that it makes pre-computations for route planning very difficult as the starting time is not known in advance.

The second approach, originally from [21], is called *time-expanded*. There, idea is to remove any time dependence from the graph by creating additional nodes for every event at a station. A node then also has a time information next to its geographic location.

Definition 12. A time expanded transit graph is a graph $G = (V, E)$ with a set of events at geographic coordinates

$$V = \{(\phi, \lambda, t) | \phi \in \left(-\frac{\pi}{2}, \frac{\pi}{2}\right), \lambda \in [-\pi, \pi), t \text{ time}\},$$

for example a train arriving or departing at a train station at a certain time.

For a node $v \in V$, v_ϕ and v_λ denote its location and v_t its time.

There is an edge $(u, w, v) \in E$ iff

1. there is a vehicle departing from u at time u_t which arrives at v at time v_t without stops in between, or
2. v is the node at the same coordinates than u with the smallest time v_t that is still greater than u_t . This edge represents exiting a vehicle and waiting for another connection. That is

$$\begin{aligned} \forall v' \in V \setminus \{v\} : v'_\phi = u_\phi \wedge v'_\lambda = u_\lambda \wedge v'_t \geq u_t \\ \Rightarrow v'_t - u_t > v_t - u_t. \end{aligned}$$

The weight w of an edge (u, w, v) is the difference between both nodes times, that is

$$w = v_t - u_t.$$

Note that weights are still positive since $v_t \geq u_t$ always holds due to construction.

Definition 12 defines such a time expanded transit graph and **Fig. 4** shows an example. For simplicity it is assumed that the trains have no stops other than shown in the schedule. The schedule lists four trains:

1. The ICE 104 which travels from Freiburg Hbf to Karlsruhe Hbf via Offenburg,
2. the RE 17024 connecting Freiburg Hbf with Offenburg,
3. the RE 17322 driving from Offenburg to Karlsruhe Hbf and
4. ICE ICE 79 which travels in the opposite direction, connecting Karlsruhe Hbf with Freiburg Hbf without intermediate stops.

As seen in the example the resulting graph has no time dependency anymore and is static, as well as all edge weights. The downside is that the graph size dramatically increases as a new node is introduced for every single event. In order to limit the growth we assume that a schedule is the same every day and does not change. In fact, most schedules are stable and often change only slightly, for example on weekends or at holidays. In practice hybrid models can be used for those exceptions.

However, the model still lacks an important feature. It does not represent *transfer buffers* [18, 4] yet. It takes some minimal amount of time to exit a vehicle and enter a different vehicle, possibly even at a different platform.

We model that by further distinguishing the nodes by arrival and departure events. In between we can then add transfer nodes which model the transfer duration. Therefore, the previous definition is adjusted and **Definition 13** is received.

Definition 13. A realistic time expanded transit graph is a graph $G = (V, E)$ with a set of events at geographic coordinates

$$V = \{(\phi, \lambda, t, e) | \phi \in \left(-\frac{\pi}{2}, \frac{\pi}{2}\right), \lambda \in [-\pi, \pi], t \text{ time}, e \in \{\text{arrival, departure, transfer}\}\},$$

→	Freiburg Hbf	Offenburg		Karlsruhe Hbf
	departure	arrival	departure	arrival
ICE 104	3:56 pm	4:28 pm	4:29 pm	4:58 pm
RE 17024	4:03 pm	4:50 pm		
RE 17322			4:35 pm	5:19 pm
←	arrival	departure	arrival	departure
ICE 79	8:10 pm			7:10 pm

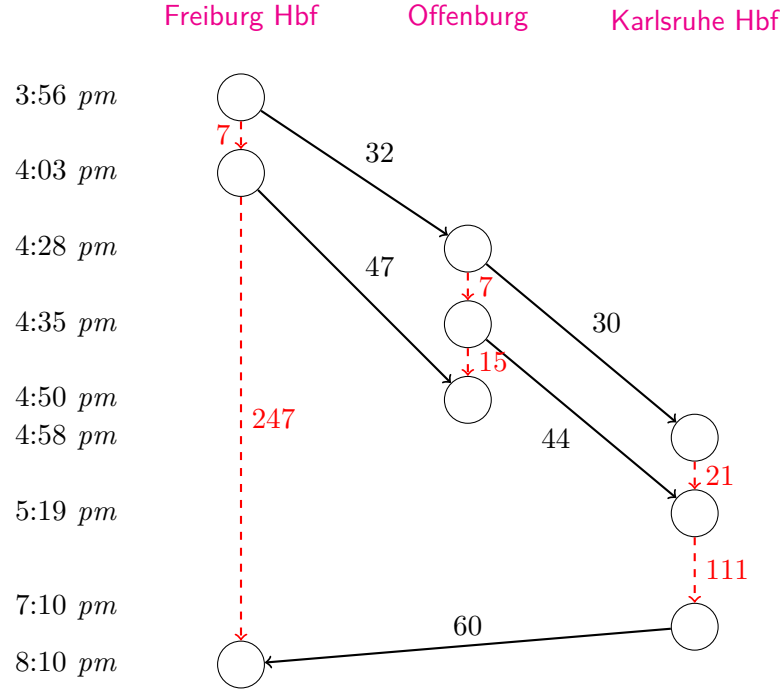


Fig. 4: Example of a transit network with its corresponding time expanded transit graph. The table shows an excerpt of a train schedule. Regular edges indicate a train connection and dashed edges waiting edges. Edge weights are measured in minutes.

for example a train arriving at a train station at a certain time.

A node $(\phi, \lambda, t, e) \in V$ is an arrival node if $e = \text{arrival}$, analogously it is a departure node for $e = \text{departure}$ and a transfer node for $e = \text{transfer}$. For a node $v \in V$, v_ϕ and v_λ denote its location, v_t its time and v_e its event type.

For every arrival node n there must exist a transfer node m at the same coordinates such that $m_t = n_t + d$ with d being the average transfer duration at the corresponding stop.

There is an edge $(u, w, v) \in E$ iff

1. $u_e = \text{departure} \wedge v_e = \text{arrival}$ such that there is a vehicle departing from u at time u_t which arrives at v at time v_t without stops in between; or
2. $u_e = \text{arrival} \wedge v_e = \text{departure}$ such that u and v belong to the same connection. For example a train arriving at a station and then departing again; or
3. $u_e = \text{arrival} \wedge v_e = \text{transfer}$ such that v is the first transfer node at the same coordinates whose time v_t comes after u_t . That is

$$\begin{aligned} \forall v' \in V \setminus \{v\} : v'_\phi = u_\phi \wedge v'_\lambda = u_\lambda \wedge v'_e = \text{transfer} \wedge v'_t \geq u_t \\ \Rightarrow v'_t - u_t > v_t - u_t. \end{aligned}$$

Such an edge represents exiting the vehicle and getting ready to enter a different vehicle; or

4. $u_e = \text{transfer} \wedge v_e = \text{transfer}$ such that v is the first transfer node at the same coordinates whose time v_t comes after u_t , representing waiting at a stop; or
5. $u_e = \text{transfer} \wedge v_e = \text{departure}$ such that u is the last transfer node at the same coordinates whose time u_t comes before v_t , i.e.

$$\begin{aligned} \forall u' \in V \setminus \{u\} : u'_\phi = v_\phi \wedge u'_\lambda = v_\lambda \wedge u'_e = \text{transfer} \wedge u'_t \leq v_t \\ \Rightarrow v_t - u'_t > v_t - u_t. \end{aligned}$$

An edge like this represents entering a different vehicle from a stop after transferring or waiting at the stop.

The weight w of an edge (u, w, v) is the difference between both nodes times, that is

$$w = v_t - u_t.$$

Fig. 5 shows how the transit graph from **Fig. 4** changes with transfer buffers.

The weight of edges connecting arrival nodes with transfer nodes is equal to the transfer duration, 5 minutes in the example. The transfer duration can be different for each edge. A transfer is now possible if the departure of the desired vehicle is after the arrival of the current vehicle plus the duration time. As seen in the example, edges connecting transfer nodes with departure nodes are present exactly in this case. A transfer from **ICE 104** to **RE 17322** in **Offenburg** is indicated by taking the edge to the first transfer node in **Offenburg** and then following the edge with cost 2 to the departure node of the train.

3.3 Link graph

In this section we examine how a road and a transit graph can be combined into a single graph such that all connections of the real network are preserved.

The approach is simple, selected nodes in the road network are connected to nodes

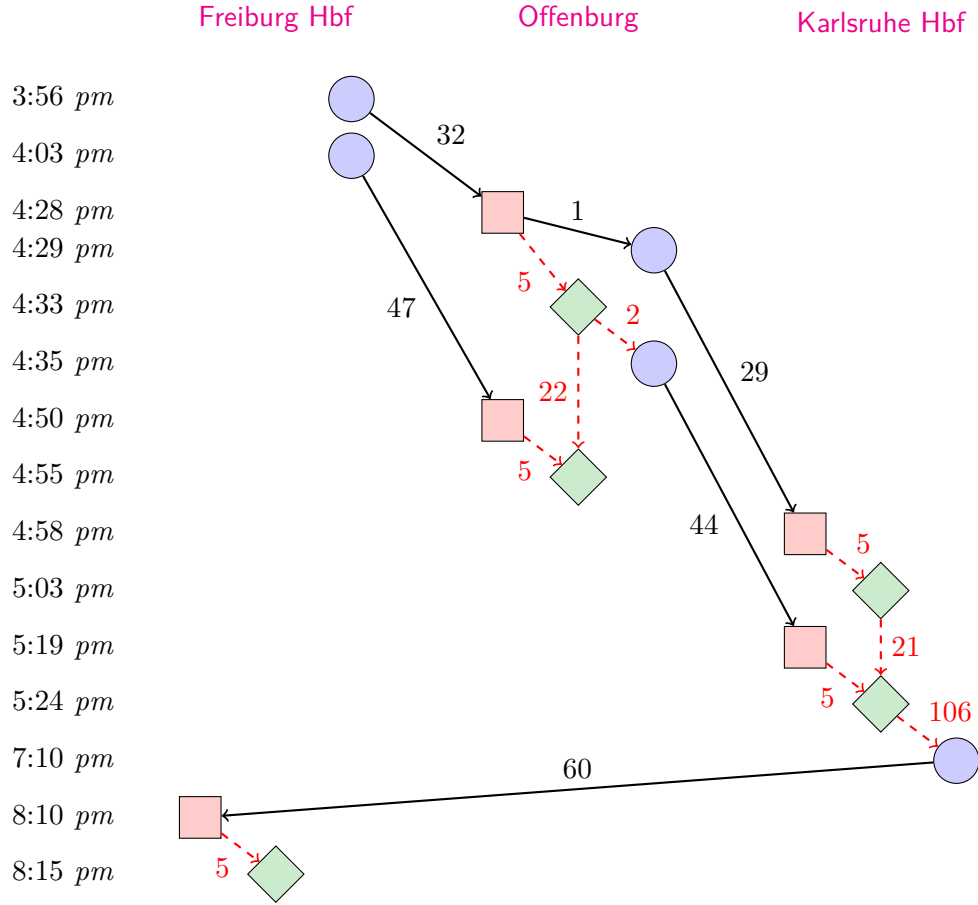


Fig. 5: Illustration of a realistic time expanded transit graph representing the schedule from **Fig. 4**. A transfer duration of 5 minutes is assumed at every stop. Rectangular nodes are arrival nodes, circular nodes represent departure nodes and diamond shaped nodes are transfer nodes. Regular edges indicate a train connection and dashed edges involve transfer nodes. Edge weights are measured in minutes.

of a certain stop in the transit network and vice versa. Since starting time is not known in advance the graph must connect a road node to all arrival nodes of a stop (compare to [9]).

In order to not miss a connection the transit graph must ensure that every connection starts with an arrival node. In **Fig. 5** this is not the case and all four trains start at a departure node. However, this is easily fixed by adding an additional arrival node to the beginning of every connection not starting with an arrival node already. The arrival nodes time is the same as the time of the departure node and both are connected by an edge with a weight of 0. **Definition 14** formalized the model.

Definition 14. Assume a road graph $R = (V_R, E_R)$, a realistic time expanded transit

graph $T = (V_T, E_T)$ where every connection in T starts by an arrival node and a partial function $\text{link} : V_R \rightarrow M$ where M contains subsets $S \subseteq V_T$. For every element $S \in M$ with an arbitrary element $s \in S$ the following properties must hold:

1. All contained elements must be arrival nodes and have the same location than s , i.e.

$$\forall s' \in S : s'_e = \text{arrival} \wedge s'_\phi = s_\phi \wedge s'_\lambda = s_\lambda.$$

2. The set must contain all arrival nodes at the location of s , i.e.

$$\nexists v \in V_T \setminus S : v_e = \text{arrival} \wedge v_\phi = s_\phi \wedge v_\lambda = s_\lambda.$$

Then a link graph is a graph $L = (V_R \cup V_T, E_R \cup E_T \cup E_L)$ with an additional set of link edges $E_L = V_R \times \mathbb{R}_{\geq 0} \times V_T$.

There is an edge $(u, 0, v) \in E_L$ iff $\text{link}(u)$ is defined and $v \in \text{link}(u)$.

The function link can be obtained in different ways. For example by creating a mapping from a road node u to a stop S if u is in the vicinity of S according to the `asTheCrowFlies` metric.

Another straightforward possibility is to always connect a stop with the road node nearest to it. We will explore this problem in **Section 4**. An obvious downside of this approach is that the nearest road node might not always have a good connectivity to the road network. A solution consists in creating a road node at the coordinates of the stop as representative. The node can then be connected with all road nodes in the vicinity.

3.4 Timetable

Timetables [4] are non-graph based representations for transit networks. They consist of stops, trips, connections and footpaths.

Definition 15. A timetable is a tuple (S, T, C, F) with stops S , trips T , connections C and footpaths F .

A stop is a position where passengers can enter or exit a vehicle, for example a train station or bus stop. It is represented as geographical coordinate (ϕ, λ) with $\phi \in (-\frac{\pi}{2}, \frac{\pi}{2})$, $\lambda \in [-\pi, \pi)$.

A trip is a scheduled vehicle, like the **ICE 104** in the example schedule of **Fig. 4** or a bus.

In contrast to a trip a connection is only a segment of a trip without stops in between. For example the connection of the **ICE 104** from **Freiburg Hbf** at 3:56 pm to **Offenburg** with arrival at 4:28 pm. It is defined as tuple $c = (s_{\text{dep}}, s_{\text{arr}}, t_{\text{dep}}, t_{\text{arr}}, o)$ with $s_{\text{dep}}, s_{\text{arr}} \in S$ representing the departure and arrival stop of the connection respectively.

Analogously t_{dep} is the time the vehicle departs at s_{dep} and t_{arr} when it arrives at s_{arr} . And $o \in T$ is the trip the connection belongs to.

Footpaths represent transfer possibilities between stops and are formalized as ordered tuple $(s_{\text{dep}}, d, s_{\text{arr}})$ with $s_{\text{dep}}, s_{\text{arr}} \in S$ the stops the footpath connects. The duration it needs to take the path by foot is represented by d , measured in seconds. Together with the set of stops S the footpaths build a graph $G = (S, F)$ representing directed edges between stops.

We require the following for the footpaths:

1. The footpaths must be transitively closed, that is

$$\exists(a, d_1, b), (b, d_2, c) \in F \Rightarrow (a, d_3, c) \in F$$

for arbitrary durations d_1, d_2, d_3 .

2. The triangle inequality must hold for all footpaths:

$$\exists(a, d_1, b), (b, d_2, c) \in F \Rightarrow \exists(a, d_3, c) \in F : d_3 \leq d_1 + d_2$$

3. Every stop must have a self-loop footpath, i.e.

$$\forall s \in S \Rightarrow (s, d, s) \in F.$$

The duration d models the transfer time at this stop, as already seen in **Section 3.2**.

The first property can easily make the set of footpaths huge. However, it is necessary for our algorithms that the amount of footpaths stays relatively small. In practice, we therefore connect each stop only to stops in its vicinity and then compute the transitive closure to ensure it is transitively closed.

To familiarize more with the model we take a look at the schedule from **Fig. 4** again. The corresponding timetable consists of:

$$S = \{f, o, k\},$$

where f, o, k represent **Freiburg Hbf**, **Offenburg** and **Karlsruhe Hbf** respectively;

$$T = \{t_{104}, t_{17024}, t_{17322}, t_{79}\},$$

representing the four trains **ICE 104**, **RE 17024**, **RE 17322** and **ICE 79**; the connections

$$\begin{aligned} &(f, o, 3:56 \text{ pm}, 4:28 \text{ pm}, t_{104}), \\ &(o, k, 4:29 \text{ pm}, 4:58 \text{ pm}, t_{104}), \\ &(f, o, 4:03 \text{ pm}, 4:50 \text{ pm}, t_{17024}), \\ &(o, k, 4:35 \text{ pm}, 5:19 \text{ pm}, t_{17322}), \\ &(k, f, 7:10 \text{ pm}, 8:10 \text{ pm}, t_{79}) \end{aligned}$$

and at least the footpaths

$$\begin{aligned} &(f, 300, f), \\ &(o, 300, o), \\ &(k, 300, k) \end{aligned}$$

for transferring at the same stop with a duration of 300 seconds (5 minutes).

If we would decide that **Offenburg** is reachable from **Freiburg Hbf** by foot and analogously **Karlsruhe Hbf** from **Offenburg**, we would also need to add a footpath connecting **Freiburg Hbf** directly with **Karlsruhe Hbf**. Else the footpaths would not be transitively closed anymore.

4 Nearest neighbor problem

In this section we introduce the **NEAREST NEIGHBOR PROBLEM**, also known as nearest neighbor search (**NNS**). First, we define the problem. Then a short overview of related research is given, after which we elaborate on a solution called **COVER TREE** [6].

Definition 16. *Given a metric space (M, d) (see **Definition 8**) with $|M| \geq 2$ and a point $x \in M$, the nearest neighbor problem asks for finding a point $y \in M$ such that*

$$y = \arg \min_{y' \in M \setminus \{x\}} d(x, y').$$

The point y is called nearest neighbor of x .

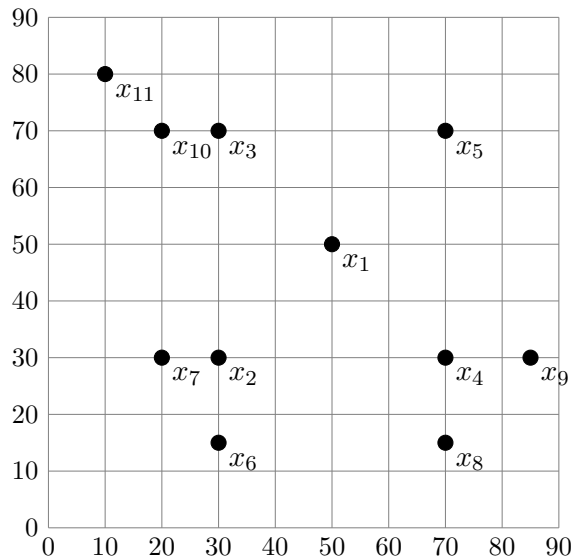


Fig. 6: Grid showing eleven points in the cartesian plane \mathbb{R}^2 .

For following examples the toy data set shown in **Fig. 6** is introduced. It consists of the points

$$\begin{aligned} x_1 &= (50, 50), \\ x_2 &= (30, 30), \\ x_3 &= (30, 70), \\ x_4 &= (70, 30), \\ x_5 &= (70, 70), \\ x_6 &= (30, 15), \\ x_7 &= (20, 30), \\ x_8 &= (70, 15), \\ x_9 &= (85, 30), \\ x_{10} &= (20, 70), \\ x_{11} &= (10, 80). \end{aligned}$$

All points are elements of the cartesian plane \mathbb{R} . The euclidean distance d is chosen as metric on this set. For two dimensions it can be defined as:

$$d : \mathbb{R}^2 \times \mathbb{R}^2, ((x_1, y_1), (x_2, y_2)) \mapsto \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Informally d computes the *ordinary* straight-line distance between two points.

The nearest neighbor of x_5 is x_1 as

$$\begin{aligned} d(x_5, x_1) &= \sqrt{(50 - 70)^2 + (50 - 70)^2} \\ &= \sqrt{800} \end{aligned}$$

is smaller than all other distances to x_5 , like

$$\begin{aligned} d(x_5, x_4) &= \sqrt{(70 - 70)^2 + (30 - 70)^2} \\ &= \sqrt{1600}. \end{aligned}$$

On the other hand, x_1 has four smallest neighbors:

$$d(x_1, x_2) = d(x_1, x_3) = d(x_1, x_4) = d(x_1, x_5)$$

Any of them is a valid solution to the nearest neighbor problem of x_1 .

The search for a nearest neighbor is a well understood problem [2, 1] and has many applications. Without restrictions, solving the problem on general metrics was proven to require $\Omega(n)$ time [2], where n is the amount of points.

Typical approaches divide the space into regions, exploiting properties of the metric

space. Common examples include **K-D TREES** [5], **VP TREES** [23], **BK-TREES** [7] and **COVER TREES** [6].

The problem also has a lot of variants. We elaborate on two of them:

Definition 17. *The k-nearest neighbors of a point $x \in M$ are the k closest points $\{y_1, y_2, \dots, y_k\} \subseteq M$ to x . That is*

$$\begin{aligned} y_1 &= \arg \min_{y' \in M \setminus \{x\}} d(x, y'), \\ y_2 &= \arg \min_{y' \in M \setminus \{x, y_1\}} d(x, y'), \\ &\vdots \\ y_k &= \arg \min_{y' \in M \setminus \{x, y_1, \dots, y_{k-1}\}} d(x, y'). \end{aligned}$$

Definition 18. *The k-neighborhood of a point $x \in M$ is the set*

$$\{y \in M \setminus \{x\} \mid d(x, y) \leq k\}.$$

4.1 Cover tree

Definition 19. *A cover tree T on a metric space (M, d) is a leveled tree (V, E) .*

The root is placed at the greatest level, denoted by $i_{\max} \in \mathbb{Z}$. The level of a node $v \in V$ is

$$\text{lvl}(v) = k - \text{depth}(v).$$

The lowest level is denoted by i_{\min} . Every node $v \in V$ is associated with a point $m \in M$. We write $\text{assoc}(v) = m$. Nodes of a certain level form a cover of points in M . A cover for a level i is defined as

$$C_i = \{m \in M \mid \exists v \in V : \text{lvl}(v) = i \wedge \text{assoc}(v) = m\}.$$

The following properties must hold

1. *For a level i there must not exist nodes which are associated with the same point $m \in M$:*

$$\nexists v, v' \in V : i = \text{lvl}(v) = \text{lvl}(v') \wedge v \neq v' \wedge \text{assoc}(v) = \text{assoc}(v')$$

So each point can at most appear once per level.

2. *$C_i \subset C_{i-1}$. This ensures that, once a point was associated with a node in a level, it appears in all lower levels too.*
3. *Points are covered by their parents:*

$$\forall p \in C_{i-1} \exists q \in C_i : d(p, q) < 2^i$$

and the node v_p with $\text{lvl}(v_p) = i \wedge \text{assoc}(v_p) = p$ is the parent of the node v_q with $\text{lvl}(v_q) = i - 1 \wedge \text{assoc}(v_q) = q$.

4. Points in a cover C_i have a separation of at least 2^i , i.e.

$$\forall p, q \in C_i : p \neq q \Rightarrow d(p, q) > 2^i.$$

A cover tree [6] has interesting distance properties on its nodes which allows for efficient retrieval of nearest neighbors. The general approach is straightforward. Given a node v in the tree placed at level i , we know that all nodes of the subtree rooted at v are associated with points inside a distance of at most 2^i . This means that, if we search for a nearest neighbor, and traverse to a node v in the tree, all nodes underneath v are relatively close to v . So, if we already have a candidate for a nearest neighbor, with distance of d and v is already further away than $d + 2^i$; v and all nodes in its subtree can not improve the distance.

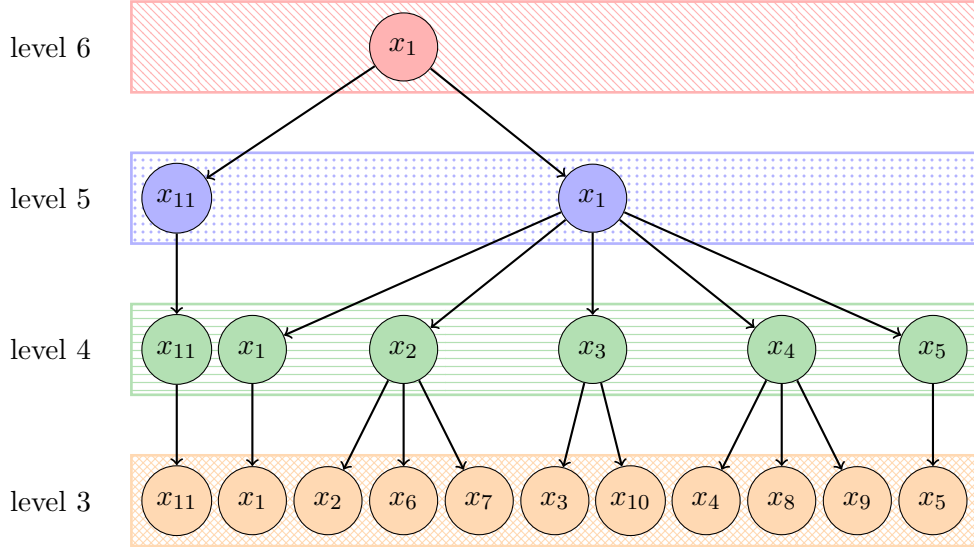


Fig. 7: Cover tree for the data set of **Fig. 6**. Nodes are vertically grouped by their levels and highlighted accordingly.

Fig. 7 shows a valid cover tree for the toy example illustrated by **Fig. 6**. The covers are

$$\begin{aligned} C_6 &= \{x_1\}, \\ C_5 &= \{x_1, x_{11}\}, \\ C_4 &= \{x_1, x_2, x_3, x_4, x_5, x_{11}\}, \\ C_3 &= \{x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9, x_{10}, x_{11}\}. \end{aligned}$$

Clearly the first property holds, there is no level where a x_i is associated with a node

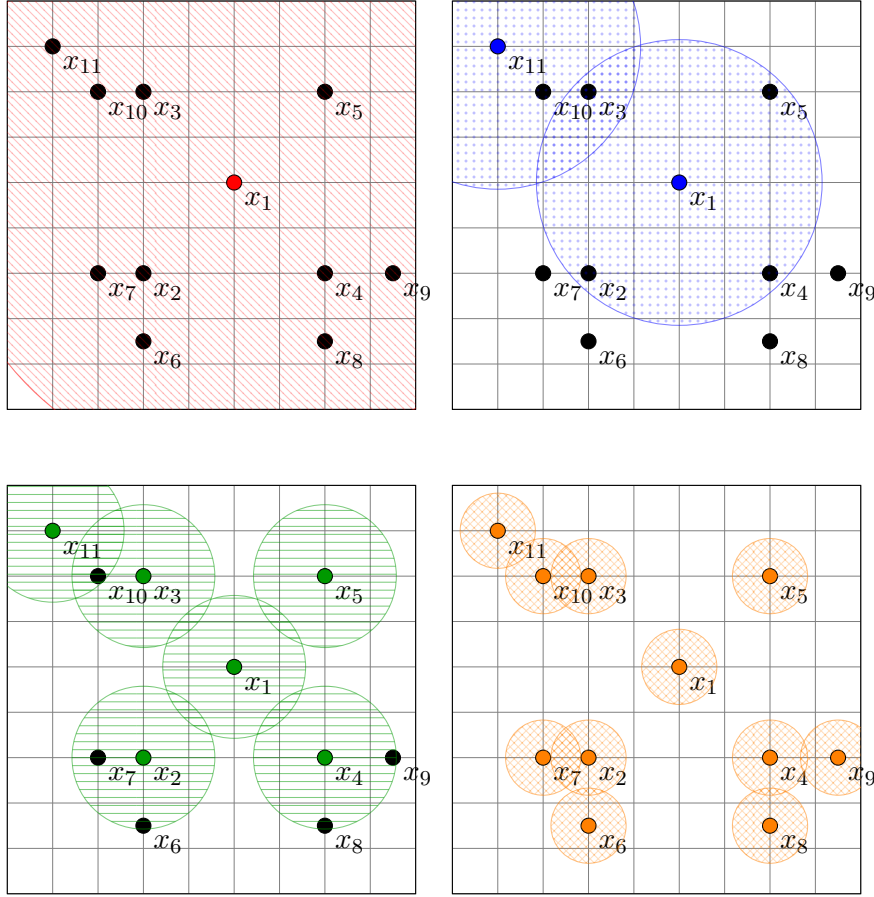


Fig. 8: Figure that shows the separation property for each level of the cover tree shown by **Fig. 7**. The levels are highlighted in the same manner than in the previous example. The levels are 6, 5, 4 and 3 from top left to bottom right. The radii around the points have a size of 2^6 , 2^5 , 2^4 and 2^3 .

more than once. The second property holds too, it is

$$C_6 \subset C_5 \subset C_4 \subset C_3.$$

For the last two properties we take a look at **Fig. 8**. It illustrates the fourth property. The property states that all points in a cover C_i must have a distance of at least 2^i to each other. For level 6 this is trivial since the set only contains x_1 . For level 5 it must hold that

$$d(x_1, x_{11}) = 50 > 32 = 2^5,$$

which is true. If this would not be the case, the figure would show the nodes included inside the circle around the other node. Analogously all nodes in C_4 and C_3 are separated enough from each other.

The third property can easily be confirmed using the figure too. It states that a node in level $i - 1$ must be closer than 2^i to its parent. Obviously this holds for x_1 and x_{11} in level 5, as a radius of 2^6 around their parent x_1 covers all nodes. Likewise are x_1, x_2, x_3, x_4 and x_5 included in the circle around their parent x_1 with radius 2^5 .

Note that it is not necessary that a node covers its whole subtree in its level. As example, we refer to x_1 in level 5 which does not cover x_{10} , as $d(x_1, x_{10}) > 2^5$, though it is part of the subtree rooted at x_1 . The third property only demands that a parent covers all its direct children, not grandchildren or similar.

Algorithm 1: Inserting a point into a cover tree operating on a metric space (M, d) .

```

input : point  $p \in M$ , candidate cover set  $Q_i \subseteq C_i$ , level  $i$ 
output: true if  $p$  was inserted at level  $i - 1$ , false otherwise

1  $Q \leftarrow \{\text{children}(q) | q \in Q_i\};$ 
2 if  $d(p, Q) > 2^i$  then
3   return false; // Check separation
4 else
5    $Q_{i-1} \leftarrow \{q \in Q | d(p, q) \leq 2^i\};$  // Covering candidates
6   if  $\neg \text{insert}(p, Q_{i-1}, i - 1) \wedge d(p, Q_i) \leq 2^i$  then
7     pick any  $q \in Q_i : d(p, q) \leq 2^i$ ;
8     append  $q$  as child to  $q$ ;
9     return true;
10  else
11    return false;

```

The cover tree is constructed using **Algorithm 1** with the maximal level i_{\max} and the cover set C_k which only consists of the root. The algorithm is stated recursively, but can easily be implemented without recursion by descending the levels and only following relevant candidates.

A point p can be appended in level $i - 1$ to a parent q in level i if the point has enough separation to all other nodes in this level, meaning more than 2^{i-1} , and is covered by the parent, that is a distance of less than 2^i . The algorithm searches such a point by descending the levels, computing the separation and appending it to a node if it also covers the point.

A search for a nearest neighbor follows a similar approach. **Algorithm 2** starts at the root and traverses the tree by following the children. The candidate set is refined by

Algorithm 2: Searching a nearest neighbor in a cover tree operating on a metric space (M, d) .

input : point $p \in M$
output: a nearest neighbor to p in M

```

1  $Q_{i_{\max}} \leftarrow C_{i_{\max}};$ 
2 for  $i$  from  $i_{\max}$  to  $i_{\min}$  do
3    $Q \leftarrow \{\text{children}(q) | q \in Q_i\};$ 
4    $Q_{i-1} \leftarrow \{q \in Q | d(p, q) \leq d(p, Q) + 2^i\};$ 
5 return  $\arg \min_{q \in Q_{i_{\min}}} d(p, q);$ 

```

only following children which are closer than

$$d(p, Q) + 2^i.$$

There, the distance to the set represents the distance of the currently best candidate. Nodes in the subtree rooted at a child can maximally be 2^i closer than the child itself. Therefore, take a look at **Fig. 8** where x_2 is maximally 2^5 closer to x_7 than x_1 , else it would not be covered by its parent x_1 . Because of that the algorithm only follows children which can have nodes in their subtree that improve over the currently best candidate. Other children are rejected.

Note that the algorithm must track down all levels, as another node could show up in the lowest level because of the separation property.

Algorithm 3: Searching the k -nearest neighbors in a cover tree operating on a metric space (M, d) .

input : point $p \in M$, amount $k \in \mathbb{N}$
output: k -nearest neighbors to p in M

```

1  $Q_{i_{\max}} \leftarrow C_{i_{\max}};$ 
2 for  $i$  from  $i_{\max}$  to  $i_{\min}$  do
3    $Q \leftarrow \{\text{children}(q) | q \in Q_i\};$ 
4   perform a  $k$ -partial sort of  $Q$ , ascending in  $d(p, q)$ ;
5   let  $q'$  be the  $k$ -th element of  $Q$ ;
6    $Q_{i-1} \leftarrow \{q \in Q | d(p, q) \leq d(p, q') + 2^i\};$ 
7 perform a  $k$ -partial sort of  $Q_{i_{\min}}$ , ascending in  $d(p, q)$ ;
8 return first  $k$  elements of  $Q_{i_{\min}};$ 

```

The cover tree can also be used to efficiently compute the k -nearest neighbors or the

Algorithm 4: Computing the k -neighborhood by using a cover tree which operates on a metric space (M, d) .

input : point $p \in M$, radius $k \in \mathbb{R}_{\geq 0}$
output: k -neighborhood of p in M

```

1  $Q_{i_{\max}} \leftarrow C_{i_{\max}};$ 
2 for  $i$  from  $i_{\max}$  to  $i_{\min}$  do
3    $Q \leftarrow \{\text{children}(q) | q \in Q_i\};$ 
4    $Q_{i-1} \leftarrow \{q \in Q | d(p, q) \leq k + 2^i\};$ 
5 return  $\{q \in Q_{i_{\min}} | d(p, q) \leq k\};$ 

```

k -neighborhood. In order to compute the k -nearest neighbors, **Algorithm 3** extends the range bound from the currently best candidate to the k -th best candidate. Likewise does **Algorithm 4** extend the bound to the given range k instead of involving candidate distances.

For other operations and a detailed analysis of the cover tree, as well as its complexity and a comparison against other techniques, refer to [6].

5 Shortest path problem

For route planning, routes through a network must be optimized in regards to one or even many criteria. A common criteria is the *travel time*. Others include cost, number of transfers or restrictions in transportation types.

In this chapter, we will first give an informal description of the **EARLIEST ARRIVAL PROBLEM**. Followed by the **SHORTEST PATH PROBLEM**, which is equivalent to the **EARLIEST ARRIVAL PROBLEM** for our graph-based network representations.

Then, we introduce algorithms for solving the problem. First, for time-independent networks, then for time-dependent. Afterwards, we explain two solutions for combined networks, using multiple transportation modes. There, the problem description slightly changes by adding transportation mode restrictions.

Definition 20. *The earliest arrival problem asks for finding a route in a network with following properties*

1. *The route must start at s and end at t .*
2. *The departure time at s is τ .*
3. *All other applicable routes must have a greater travel time, i.e. arrive later at t .*

Points s and t are given source and target points in the network respectively. τ is the desired departure time, it may be ignored for a time-independent network.

Definition 21. Given a graph $G = (V, E)$, source and target nodes $s, t \in V$ and a desired departure time τ , the shortest path problem asks for a path p (see **Definition 3**) which

1. begins at s and ends at t ,
2. has the smallest weight of all applicable paths.

The arrival time at t is τ plus the weight of p . In a time-dependent graph τ must be used to ensure correct edge weights. The path p is called shortest path.

Additionally, we consider a special variant of the shortest path problem:

Definition 22. The many-to-one shortest path problem is a variation of the shortest path problem where the source consists of a set of source nodes $S \subseteq V$.

The problem asks for the path p that starts at the source $s \in S$ which minimizes the path weight.

5.1 Time-independent

Route planning in time-independent networks is a very well studied problem. Many efficient solutions to the shortest path problem exists. We introduce a very basic algorithm, **DIJKSTRA** and a simple improvement based on heuristics, **A***.

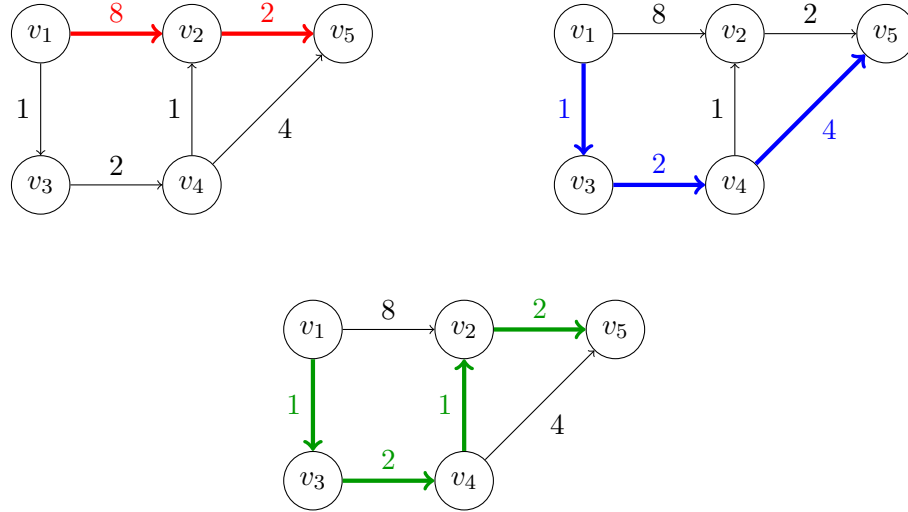


Fig. 9: Example for a time independent network, represented by a road graph. The figure shows three paths from v_1 to v_5 . From top left to bottom right, the path weights are 10, 7 and 6. The last example represents the shortest path from v_1 to v_5 .

The network shown by **Fig. 9** acts as toy example for this section.

5.1.1 Dijkstra

DIJKSTRA [8] is a simple approach to solving the shortest path problem. It can be viewed as the logical extension of breadth-first search (**BFS**) [8] in weighted graphs. The algorithm revolves around a priority queue where it stores neighboring nodes, sorted by their shortest path cost. In each round, the node with the smallest shortest path cost is *relaxed*. That is, all its neighboring, not already relaxed, nodes are added to the queue. The algorithm terminates as soon as the target node has been relaxed. **Algorithm 5** gives a formal description.

To familiarize with the algorithm, we step through the execution for the graph shown by **Fig. 9**, with v_1 as source and v_5 as target node.

The **dist** function, often implemented as array, stores the tentative shortest path weight to the given node. **prev** is used for path extraction at the end, it stores the parent nodes used for the shortest paths represented by **dist**. The algorithm starts by initializing both collections with default values. Initially the distance to all nodes, except the source, is unknown. Thus, ∞ is used for them. Q represents the list of nodes that need to be processed, usually implemented as priority queue. Initially, it only holds the source node s .

In the example Q is initially $\{v_1\}$. The algorithm then relaxes v_1 and stores distances to its neighbors:

$$\begin{aligned} \text{dist}(v_2) &= 8 & \text{prev}(v_2) &= v_1, \\ \text{dist}(v_3) &= 1 & \text{prev}(v_3) &= v_1 \end{aligned}$$

Additionally, the queue Q is updated, it is

$$Q = \{v_2, v_3\}.$$

The next iteration of the loop starts and the node with the smallest distance is chosen, i.e. v_3 . The node is relaxed and we receive

$$\begin{aligned} \text{dist}(v_4) &= 3 & \text{prev}(v_4) &= v_3, \\ Q &= \{v_2, v_4\}. \end{aligned}$$

The next node is v_4 , yielding

$$\begin{aligned} \text{dist}(v_2) &= 4 & \text{prev}(v_2) &= v_4, \\ \text{dist}(v_5) &= 7 & \text{prev}(v_5) &= v_4, \\ Q &= \{v_2, v_5\}. \end{aligned}$$

Note that v_4 improves the distance to v_2 . The previous values for v_2 are overwritten and the tentative shortest path to v_2 uses $(v_4, 1, v_2)$ and not $(v_1, 8, v_2)$ anymore. In the next round v_2 is relaxed which improves the distance to v_5 :

$$\begin{aligned} \text{dist}(v_5) &= 6 & \text{prev}(v_5) &= v_2, \\ Q &= \{v_5\}. \end{aligned}$$

Algorithm 5: Dijkstra's algorithm for computing shortest paths in time-independent graphs.

```

input : graph  $G = (V, E)$ , source  $s \in V$ , target  $t \in V$ 
output: shortest path from  $s$  to  $t$ 

// Initialization
1 for  $v \in V$  do
2    $\text{dist}(v) \leftarrow \infty$ ;
3    $\text{prev}(v) \leftarrow \text{undefined}$ ;

4  $\text{dist}(s) \leftarrow 0$ ;
5  $Q \leftarrow \{s\}$ ;

// Compute shortest paths
6 while  $Q$  is not empty do
7    $u \leftarrow \arg \min_{u' \in Q} \text{dist}(u')$ ;
8    $Q \leftarrow Q \setminus \{u\}$ ;
9   if  $u == t$  then
10    break;

    // Relax  $u$ 
11   for outgoing edge  $(u, w, v) \in E$  do
12      $\text{currentDist} \leftarrow \text{dist}(u) + w$ ;
13     if  $\text{currentDist} < \text{dist}(v)$  then
14       // Improve distance by using this edge
15        $\text{dist}(v) \leftarrow \text{currentDist}$ ;
16        $\text{prev}(v) \leftarrow u$ ;
17        $Q \leftarrow Q \cup \{v\}$ ;

// Extract path by backtracking
17  $p \leftarrow \text{empty path}$ ;
18  $u \leftarrow t$ ;
19 while  $\text{prev}(u) \neq \text{undefined}$  do
20    $w \leftarrow \text{dist}(u) - \text{dist}(\text{prev}(u))$ ;
21   prepend  $(\text{prev}(u), w, u)$  to  $p$ ;
22    $u \leftarrow \text{prev}(u)$ ;
23 prepend  $s$  to  $p$ ;
24 return  $p$ ;
```

The only node left is the target node v_5 now. It is relaxed and the loop terminates. The algorithm backtracks the parent pointer

```

prev( $v_5$ ) =  $v_2$ ,
prev( $v_2$ ) =  $v_4$ ,
prev( $v_4$ ) =  $v_3$ ,
prev( $v_3$ ) =  $v_1$ ,
prev( $v_1$ ) = undefined

```

and constructs the shortest path

$$p = (v_1, 1, v_3)(v_3, 2, v_4)(v_4, 1, v_2)(v_2, 2, v_5)$$

which is the path shown by the last example in the figure.

5.1.2 A* and ALT

An important observation of **DIJKSTRA** is that, if it settles the shortest path distance to a node, then, all nodes which are closer to the source, were already settled in a previous round.

Moreover, the algorithm explores the graph in all directions equally. It has no sense of *goal direction*.

The **A*** algorithm [15] is a simple extension of **DIJKSTRA** which improves its efficiency by steering the exploration more towards the target. **Fig. 10** illustrates this by comparing the *search space* of both algorithms. The search space of **A*** is smaller and much more directed to the target node t .

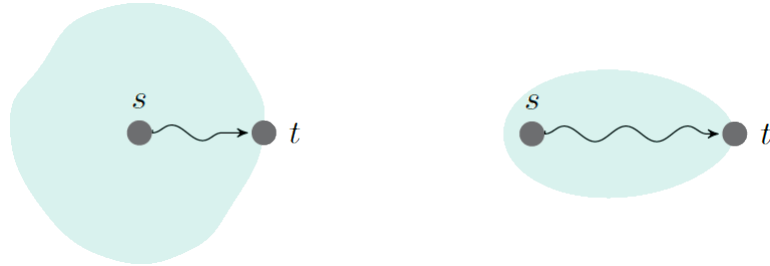


Fig. 10: Schematic illustration of a query processed by **DIJKSTRA** (left) and **A*** (right). The highlighted areas indicate the *search space*, i.e. the nodes the algorithm has explored already. The illustration is from [4].

Unfortunately, computing the exact goal direction is as hard as computing the shortest path to the target. Therefore, a heuristic is used to approximate the direction. The choice of the heuristic heavily depends on the underlying network. In the worst case, a heuristic may not improve over **DIJKSTRA** and the same search space is received. In the best case, the algorithm explores only the nodes on the shortest path.

Such a heuristic must fulfill two properties, formulated by **Definition 23**.

Definition 23. Given a graph $G = (V, E)$, a metric dist on V (see **Definition 7**), a heuristic is a function $h : V \times V \rightarrow \mathbb{R}_{\geq 0}$ which approximates dist . The heuristic h must be

1. admissible, i.e. never overestimate:

$$\forall u, t \in V : h(u, t) \leq \text{dist}(u, t)$$

2. monotone, i.e. satisfy the triangle inequality:

$$\forall t \in V \forall (u, w, v) \in E : h(u, t) \leq w + h(v, t)$$

Given such a heuristic h , the A^* algorithm is received by adjusting **line 7** of **Algorithm 5** to

$$u \leftarrow \arg \min_{u' \in Q} \text{dist}(u') + h(u', t).$$

This will prefer nodes that are estimated to be closer to the target before others. By that, the algorithms search space first expands into a direction that minimizes the distance according to the heuristic h .

A common choice for a simple heuristic is the *as-the-crow-flies* metric (see **Definition 10**). The properties are easily verified. A theoretically shortest path has the shortest possible distance and uses the fastest available transportation mode. This is exactly the path represented by the *straight-line* distance, computed by the *as-the-crow-flies* metric. It can thus never overestimate. It is also trivially monotone since it is a metric, i.e. the triangle inequality holds for all elements.

A heuristic is a good choice if it approximates the actual shortest path distance well. As such, the *as-the-crow-flies* heuristic works well on networks with a high connectivity in all directions. For example a residential area of a city without one way streets. Unfortunately, in road networks, the common case is to first drive into the opposite direction in order to reach a fast highway. This even gets worse on networks where the importance of nodes heavily differ, such as public transit networks. For train networks, the typical case is that one first needs to travel to a main station. This is obviously due to a main station having a much better connectivity and faster trains available. Because of that, the effectiveness of *as-the-crow-flies* is very limited on such networks.

The *landmark heuristic* partially solves the issue. An A^* algorithm using the landmark heuristic is called **ALT** [15], which stands for *landmarks and triangle inequality*.

The heuristic provides a more generic approach by approximating the distance between nodes u and v by using pre-computed distances with pre-determined nodes l , called *landmarks*.

Definition 24. Given a set of landmarks $L \subseteq V$, the heuristic landmarks is defined by

$$\text{landmarks}(u, v) = \max_{l \in L} (\max\{\text{dist}(u, l) - \text{dist}(v, l), \text{dist}(l, v) - \text{dist}(l, u)\}).$$

Obviously, the heuristic improves if the set of landmarks is increased. However, actual shortest path distances from all landmarks to all other nodes in the graph must be pre-computed. With an increasing amount of landmarks the pre-computation might not be feasible anymore because it takes too long or consumes too much space. Note that if $L = V$, the heuristic becomes the actual shortest path distance function, i.e. $\text{landmarks} = \text{dist}$.

In practice, an amount between 20 and 50 randomly chosen nodes seems to be a good compromise. Refer to [15] for a detailed analysis.

The computation of the actual shortest path distances, to and from the landmarks, can be done by using **DIJKSTRA**. But, instead of running the algorithm for all pairs of nodes, the distances can be obtained with two runs only. Therefore, the algorithm is slightly modified by dropping **lines 9** and **10**, such that the algorithm relaxes the whole network. By that, a single run of **DIJKSTRA** with a landmark l as source, computes the distances $\text{dist}(l, v)$ to all nodes v in the network. By reversing the graph, i.e. edges (u, w, v) become (v, w, u) , the distances to the landmarks can be obtained analogously with l as source again. Depending on the graph implementation, reversal can be done in $\mathcal{O}(1)$ by only implicitly reversing the edges.

5.2 Time-dependent

Approaches designed for time-independent networks, such as **ALT**, have an important drawback. Optimization is always done on assuming that edge costs are constant. However, in a time-dependent network, this is not the case. The weight of an edge is dependent on the departure time, which is not known in advance.

DIJKSTRA and its variants **A*** and **ALT** can easily be adapted to also work in time-dependent networks by taking the departure time into consideration when computing the weight of an edge. However, their effectiveness is very limited. Nonetheless, they were used for a long time for time-dependent networks too. With increasing research on route planning in time-dependent networks, more effective algorithms, such as **TRANSFER PATTERNS** [3] and **CSA** [11], were developed. Many of them do not use graphs and prefer data-structures that are designed for time-dependent data, such as *timetables* (see **Section 3.4**).

5.2.1 Connection scan

Connection scan (**CSA**) [11] is an algorithm for route planning specially designed for time-dependent networks, such as public transit networks. It processes the network represented as timetable, as defined by **Definition 15**.

The algorithm is very simple. All connections of the timetable are sorted by their departure time. Give a query connections are explored increasing in their departure time. The algorithm is fast primarily due to the fact that connections can be maintained in a simple array. In contrast to **DIJKSTRA**, it does not need to maintain a priority queue or other more complex data-structures. Arrays are heavily optimized and benefit from a lot of effects, like cache locality [16].

Algorithm 6 shows the full connection scan algorithm. The array S stores for each stop the currently best arrival time. T associates for each trip the first connection it is taken with. J is used for path extraction and memorizes for each stop a segment of a trip, consisting of enter and exit connections c_{enter} and c_{exit} respectively, and a footpath f :

$$(c_{\text{enter}}, c_{\text{exit}}, f)$$

It represents a path which takes the segment of the trip starting at c_{enter} , ending at c_{exit} and then taking the footpath f from the arrival stop of c_{exit} . Such an entry is associated to the arrival stop of the footpath f , always representing the parent path that results in the current best arrival time for the corresponding stop.

The algorithm starts by initializing the arrays with default values and relaxing all initial footpaths. Connections are then explored increasing in their departure time, starting from the first connection c_0 that starts after the departure time τ . **Line 7** is typically implemented as *binary search* [17] on a sorted array of connections C .

Line 9 is the stopping criteria, which lets the algorithm terminate once a connection departs after the current best arrival time at the target t . Since connections are explored increasing in time, it is impossible that a connection can improve on the arrival time anymore.

Line 11 will only explore a connection if a previous connection of the same trip was already used, indicating traveling without a transfer; or if it was already possible to arrive at the stop earlier with a previous connection, indicating a transfer at this stop.

A connection is then only relaxed if it improves the arrival time at its arrival stop, represented by **line 14**. If so, all outgoing footpaths are explored. A footpath represents exiting the vehicle, walking to the arrival stop of the footpath ready for entering another vehicle. Note that self-loop footpaths must be contained in timetables (compare to **Definition 15**), making it possible to transfer at one stop.

Line 16 only considers footpaths that improve the arrival time at the corresponding stop. **Line 18** stores the path represented by taking this connection and the footpath.

For an example, we refer to the schedule of **Fig. 4** again. The corresponding timetable is explained in **Section 3.4**, we use the same notion again. It consists of five connections, denoted by c_1, c_2, c_3, c_4 and c_5 , sorted by departure time. We assume only the three self-loop footpaths on the stops f , o and k .

Assume a query from **Freiburg Hbf**, represented by stop f to **Karlsruhe Hbf**, represented

Algorithm 6: Connection scan algorithm for computing shortest paths in time-dependent networks, represented by timetables.

input : timetable (S, T, C, F) , source $s \in S$, target $t \in S$, departure time τ
output: shortest path from s to t

// Initialization

```

1 for  $u \in S$  do  $S[u] \leftarrow \infty$ ;
2 for  $o \in T$  do  $T[o] \leftarrow \text{undefined}$ ;
3 for  $u \in S$  do  $J[u] \leftarrow (\text{undefined}, \text{undefined}, \text{undefined})$ ;
4 for  $f = (u_{\text{dep}}, d, u_{\text{arr}}) \in F : u_{\text{dep}} = s$  do
5    $S[u_{\text{arr}}] \leftarrow \tau + d$ ;
6    $J[u_{\text{arr}}] \leftarrow (\text{undefined}, \text{undefined}, f)$ ;

// Explore connections increasing in departure time
7  $c_0 \leftarrow \arg \min_{(u_{\text{dep}}, u_{\text{arr}}, \tau_{\text{dep}}, \tau_{\text{arr}}, o) \in C : \tau_{\text{dep}} \geq \tau} \tau_{\text{dep}}$ ;
8 for  $c = (u_{\text{dep}}, u_{\text{arr}}, \tau_{\text{dep}}, \tau_{\text{arr}}, o) \in C$  increasing by  $\tau_{\text{dep}}$ , starting from  $c_0$  do
9   if  $\tau_{\text{dep}} \geq S[t]$  then
10    break;
11   if  $T[o] \neq \text{undefined} \vee \tau_{\text{dep}} \geq S[u_{\text{dep}}]$  then
12     if  $T[o] == \text{undefined}$  then
13        $T[o] \leftarrow c$ ;
14     if  $\tau_{\text{arr}} < S[u_{\text{arr}}]$  then
15       for  $f = (v_{\text{dep}}, d, v_{\text{arr}}) \in F : v_{\text{dep}} = u_{\text{arr}}$  do
16         if  $\tau_{\text{arr}} + d < S[v_{\text{arr}}]$  then
17            $S[v_{\text{arr}}] \leftarrow \tau_{\text{arr}} + d$ ;
18            $J[v_{\text{arr}}] \leftarrow (T[o], c, f)$ ;

// Extract path by backtracking
19  $p \leftarrow \text{empty path}$ ;
20  $u \leftarrow t$ ;
21 while  $c_{\text{enter}} \neq \text{undefined} : (c_{\text{enter}}, c_{\text{exit}}, f) = J[u]$  do
22   prepend  $f$  to  $p$ ;
23   prepend the part of the trip between  $c_{\text{enter}}$  and  $c_{\text{exit}}$  to  $p$ ;
24    $u \leftarrow v_{\text{dep}} : (v_{\text{dep}}, v_{\text{arr}}, \tau'_{\text{dep}}, \tau'_{\text{arr}}, o) = c_{\text{enter}}$ ;
25 prepend  $f : (\text{undefined}, \text{undefined}, f) = J[s]$  to  $p$ ;
26 return  $p$ ;
```

by k , with a departure time of $\tau = 3:50$ *pm*. The initial configuration after **line 3** is

$$\begin{aligned} S[f] &= S[o] = S[k] = \infty, \\ T[t_{104}] &= T[t_{17024}] = T[t_{17322}] = T[t_{79}] = \text{undefined}, \\ J[f] &= J[o] = J[k] = (\text{undefined}, \text{undefined}, \text{undefined}). \end{aligned}$$

Then the footpath $(f, 300, f)$ departing at **Freiburg Hbf** is relaxed, resulting in

$$\begin{aligned} S[f] &= 3:55 \text{ pm}, \\ J[f] &= (\text{undefined}, \text{undefined}, (f, 300, f)). \end{aligned}$$

Connections are now explored increasing in departure time, starting with

$$c_1 = (f, o, 3:56 \text{ pm}, 4:28 \text{ pm}, t_{104}).$$

The connection is considered since we already arrived at **Freiburg Hbf** before 3:56 *pm*. The trip is set and the footpath at **Offenburg** is relaxed, yielding

$$\begin{aligned} T[t_{104}] &= c_1, \\ S[o] &= 4:33 \text{ pm}, \\ J[o] &= (c_1, c_1, (o, 300, o)). \end{aligned}$$

The next connection is

$$c_2 = (f, o, 4:03 \text{ pm}, 4:50 \text{ pm}, t_{17024}).$$

However, it induces no changes, as the previous connection already arrived at **Offenburg** earlier. The algorithm continues by exploring

$$c_3 = (o, k, 4:29 \text{ pm}, 4:58 \text{ pm}, t_{104}).$$

The connection is considered because the trip t_{104} was used before already, indicating that the trip can be taken without transferring. Else it would not be applicable, since the current best arrival time at **Offenburg**, including the transfer duration of 5 minutes, is 4:33 *pm*, which is after the departure time of c_3 . The changes are

$$\begin{aligned} S[k] &= 5:03 \text{ pm}, \\ J[k] &= (c_1, c_3, (k, 300, k)). \end{aligned}$$

In the next iteration

$$c_4 = (o, k, 4:35 \text{ pm}, 5:19 \text{ pm}, t_{17322})$$

is considered, again inducing no changes. The algorithm then terminates exploration since the last connection

$$c_5 = (k, f, 7:10 \text{ pm}, 8:10 \text{ pm}, t_{79})$$

departs after the current best arrival time at **Karlsruhe Hbf**, which is $S[k] = 5:03 \text{ pm}$.

Path construction is straightforward, it is

$$\begin{aligned} J[k] &= (c_1, c_3, (k, 300, k)), \\ J[f] &= (\text{undefined}, \text{undefined}, (f, 300, f)), \end{aligned}$$

which yields the path which takes

the footpath from **Freiburg Hbf** to **Freiburg Hbf**,

t_{104} starting with c_1 to c_3 , which is using the **ICE 104** from **Freiburg Hbf** to **Karlsruhe Hbf**,

and a final footpath from **Karlsruhe Hbf** to **Karlsruhe Hbf**.

The earliest arrival time at **Karlsruhe Hbf** is $S[k] = 5:03 \text{ pm}$.

5.3 Multi-modal

Blabla

5.3.1 Modified Dijkstra

Blabla

5.3.2 Access nodes

Blabla

5.4 Other algorithms

Blabla

6 Evaluation

Blabla

6.1 Input data

Blabla

6.2 Experiments

Blabla

6.2.1 Nearest neighbor computation

Blabla

6.2.2 Uni-modal routing

Blabla

6.2.3 Multi-modal routing

Blabla

6.3 Summary

Blabla

7 Conclusion

Blabla

References

- [1] Mohammad Reza Abbasifard, Bijan Ghahremani, and Hassan Naderi. A survey on nearest neighbor search methods. 2014.
- [2] Alexandr Andoni. Nearest neighbor search : the old , the new , and the impossible by alexandr andoni. 2009.
- [3] Hannah Bast, Erik Carlsson, Arno Eigenwillig, Robert Geisberger, Chris Harrelson, Veselin Raychev, and Fabien Viger. Fast routing in very large public transportation networks using transfer patterns. In Mark de Berg and Ulrich Meyer, editors, *Algorithms – ESA 2010*, pages 290–301, Berlin, Heidelberg, 2010. Springer Berlin Heidelberg.
- [4] Hannah Bast, Daniel Delling, Andrew Goldberg, Matthias Müller-Hannemann, Thomas Pajor, Peter Sanders, Dorothea Wagner, and Renato F. Werneck. *Route Planning in Transportation Networks*, pages 19–80. Springer International Publishing, Cham, 2016.
- [5] Jon Louis Bentley. Multidimensional binary search trees used for associative searching. *Commun. ACM*, 18(9):509–517, September 1975.
- [6] Alina Beygelzimer, Sham Kakade, and John Langford. Cover trees for nearest neighbor. In *Proceedings of the 23rd International Conference on Machine Learning*, ICML '06, pages 97–104, New York, NY, USA, 2006. ACM.
- [7] W. A. Burkhard and R. M. Keller. Some approaches to best-match file searching. *Commun. ACM*, 16(4):230–236, April 1973.
- [8] Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to Algorithms, Third Edition*. The MIT Press, 3rd edition, 2009.

- [9] Daniel Delling, Thomas Pajor, and Dorothea Wagner. Accelerating multi-modal route planning by access-nodes. In Amos Fiat and Peter Sanders, editors, *Algorithms - ESA 2009*, pages 587–598, Berlin, Heidelberg, 2009. Springer Berlin Heidelberg.
- [10] Daniel Delling, Peter Sanders, Dominik Schultes, and Dorothea Wagner. *Engineering Route Planning Algorithms*, pages 117–139. Springer Berlin Heidelberg, Berlin, Heidelberg, 2009.
- [11] Julian Dibbelt, Thomas Pajor, Ben Strasser, and Dorothea Wagner. Connection scan algorithm. *CoRR*, abs/1703.05997, 2017.
- [12] Reinhard Diestel. *Graph Theory, 4th Edition*, volume 173 of *Graduate texts in mathematics*. Springer, 2012.
- [13] Wei Dong. An overview of in-vehicle route guidance system. In *Australasian Transport Research Forum*, volume 2011. Citeseer, 2011.
- [14] R. L. French. Historical overview of automobile navigation technology. In *36th IEEE Vehicular Technology Conference*, volume 36, pages 350–358, May 1986.
- [15] Andrew V. Goldberg and Chris Harrelson. Computing the shortest path: A search meets graph theory. In *Proceedings of the Sixteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '05, pages 156–165, Philadelphia, PA, USA, 2005. Society for Industrial and Applied Mathematics.
- [16] Dirk Grunwald, Benjamin Zorn, and Robert Henderson. Improving the cache locality of memory allocation. In *ACM SIGPLAN Notices*, volume 28, pages 177–186. ACM, 1993.
- [17] Donald E. Knuth. *The Art of Computer Programming, Volume 3: (2Nd Ed.) Sorting and Searching*. Addison Wesley Longman Publishing Co., Inc., Redwood City, CA, USA, 1998.
- [18] Matthias Müller-Hannemann, Frank Schulz, Dorothea Wagner, and Christos Zaroliagis. Timetable information: Models and algorithms. In Frank Geraets, Leo Kroon, Anita Schoebel, Dorothea Wagner, and Christos D. Zaroliagis, editors, *Algorithmic Methods for Railway Optimization*, pages 67–90, Berlin, Heidelberg, 2007. Springer Berlin Heidelberg.
- [19] John P. Snyder. Flattening the earth: Two thousand years of map projections. 10 1994.
- [20] C. C. Robusto. The cosine-haversine formula. *The American Mathematical Monthly*, 64(1):38–40, 1957.
- [21] Frank Schulz, Dorothea Wagner, and Karsten Weihe. Dijkstra’s algorithm on-line: An empirical case study from public railroad transport. *J. Exp. Algorithmics*, 5, December 2000.

- [22] Daniel Tischner. Cobweb. <https://github.com/ZabuzaW/Cobweb>, 2018.
- [23] Peter N. Yianilos. Data structures and algorithms for nearest neighbor search in general metric spaces. In *Proceedings of the Fifth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, 1993.
- [24] Yilin Zhao. *Vehicle location and navigation systems*. 1997.