

Landmark indexing for scalable evaluation of label-constrained reachability queries

L.D.J. (Lucien) Valstar

Supervisors:
G.H.L. Fletcher
Y. Yoshida

Committee members:
A. Driemel
G.H.L. Fletcher
M.A. Westenberg
Y. Yoshida

version 1.0

Eindhoven, July 2016

Abstract

Our world today is generating huge amounts of graph data such as social networks, biological networks, and the semantic web. Many of these real-world graphs are edge-labelled graphs, i.e. each edge has a label that denotes the relationship between the two vertices connected by the edge. A fundamental research problem on these is how to handle reachability in these kinds of graphs: *can vertex u reach vertex v using only edges with particular labels?* There has not been very much research on this topic yet. Hence we have come up with our own solution based on an index. We dealt with the problem by looking for several algorithms that can deal with these queries. We ran several experiments to examine their performance relative to a baseline algorithm. The results show that there is a clear performance improvement by building the index we develop here.

Keywords: LCR, Label-constrained reachability, Labelled graph, Landmarks, Index

Contents

Contents	v
1 Introduction	1
1.1 Problem motivation	1
1.2 State of the art	1
1.3 Contributions	2
1.4 Thesis outline	2
2 Problem statement	3
2.1 Reachability in graphs	3
2.2 Label-constrained reachability	4
2.3 The problem	5
2.4 Auxiliary definitions	5
3 Literature analysis	7
3.1 Reachability	7
3.1.1 2-hop cover	7
3.2 Label-constrained reachability	9
3.2.1 Bonchi et al.	9
3.2.2 Zou et al.	9
3.2.3 Fletcher and Yoshida	16
4 Methods	19
4.1 Existing methods	19
4.1.1 BFS	19
4.1.2 Best effort Zou	20
4.2 Our contributions	21
4.2.1 General comments	21
4.2.2 LandmarkedIndex	21
4.2.3 Partial	29
4.2.4 NeighbourExchange	29
4.2.5 Joindex	29
4.2.6 ClusteredExact	30
4.3 Implementation details	32
4.4 Index maintenance	32
4.4.1 Adding an edge	33
4.4.2 Removing an edge	33
4.4.3 Changing edge label	34
4.4.4 Adding a node	34
4.4.5 Removing a node	34
4.5 Extensions	34
4.5.1 Query for all nodes	34
4.5.2 Distance queries	35

5	Experimental design	37
5.1	Datasets	37
5.1.1	Synthetic datasets	37
5.1.2	Real datasets	37
5.1.3	Summary of datasets	39
5.2	Queries	41
5.3	Hardware	41
5.4	Methods	41
6	Experiments	43
6.1	Introduction	43
6.2	Part 1: small graphs ($0 < E \leq 5,000$)	44
6.2.1	Datasets and methods	44
6.2.2	Index construction time (s) and size (MB)	44
6.2.3	Total speed-ups	45
6.3	Part 2: medium graphs ($5,000 < E \leq 500,000$)	48
6.3.1	Datasets and methods	48
6.3.2	Index construction time (s)	48
6.3.3	Index size (MB)	49
6.3.4	Index construction time (s) and index size (MB) when $ \mathcal{L} \geq 8$	50
6.3.5	Total speed-ups	50
6.3.6	Total speed-ups when $ \mathcal{L} \geq 8$	52
6.3.7	Individual speed-ups	52
6.4	Part 3: large graphs ($ E > 500,000$)	56
6.4.1	LI+OTH+EXTv1: Index construction time (s) and size (MB)	56
6.4.2	LI+OTH+EXTv1: Total speed-ups	58
6.4.3	LI+OTH+EXTv2: Index construction time (s) and size (MB)	60
6.4.4	LI+OTH+EXTv2: Total speed-ups	62
6.5	Maintenance	64
6.5.1	Adding an edge	64
6.6	Extensions	64
6.6.1	Query for all nodes	64
6.6.2	Distance queries	65
7	Conclusion	67
7.1	Future work	67
	Bibliography	69
	Appendix	69

Chapter 1

Introduction

In this chapter we introduce the problem studied in this thesis, we provide an overview of the state of the art, explain briefly our contributions and give a thesis outline.

1.1 Problem motivation

Our world today is generating huge amounts of graph data such as social networks, biological networks, and the semantic web. The sheer size of these graphs (e.g. Google Knowledge Graph has 570 million vertices/objects and 18 billion edges/facts) turns a simple query into a difficult one. One of these queries is “reachability”, i.e. the question on whether we can from a point (or vertex) A reach another point B in the graph using the edges of the graph.

Many real-world graphs are edge-labelled graphs. This means that edges in the graph have a label from a pre-defined label set. This changes a reachability-query into a so-called “label-constrained reachability” query or “LCR”-query. The question for such a query is: “can we reach from point A another point B in the graph using only certain types of edges?”.

Examples of graphs for which there are labels are numerous.

In **Social networks** each person is represented as a vertex. Any person can have an interaction with another person if they are related in some way. Examples of such relationships could be colleague, friend or family (sister-of, brother-of, son-of). It can be interesting to look at more complicated relationships between people. When we want to know whether person A is a remote relative of person B we might look to explore the graph using only edges brother-of, sister-of, son-of, etc.

In **Bioinformatics** there is a need to understand metabolic chain reactions in cellular systems. For this a metabolic network is used. In such a network, each vertex represents a chemical compound. An edge indicates that one compound can transform into another. The edge label records the enzymes that are needed to control the reaction. The basic question is to whether there is a pathway between two compounds where the enzymes are present under certain conditions.

RDF-graphs consist of triplets (or facts) that indicate that entity A is related to entity B through a certain relationship. Information from multiple sources (e.g. IMDb, Facebook or Wikipedia) may be combined. A given person A might create a FOAF-record describing personal information, interests and friendships. This may be linked to for instance a movie review this user has written on another website. A query could be “how is a respondent B to the movie review written by person A related to person A ” or “in what ways are two particular movie reviews X and Y written by user A related”.

1.2 State of the art

There is not much work on the topic of “LCR”. The baseline solution is just running BFS on the graph. BFS is the baseline, as it builds no index and has the maximal query answering time. Any other method builds an index and uses this to speed up its query answering time.

Zou et al. [11] have proposed their solution. We were unable to obtain a version of their code and hence we made our own best-effort implementation of it. The time to construct their index is very low compared to those of other papers and to the index construction times of our methods.

Then there is Bonchi et al. [1]. They aimed at answering a query similar but more difficult than the query we study. Nevertheless the results in their paper can be used as a comparison for our results.

1.3 Contributions

This study contributes to a few novel things. The first is exploring the existing (ZOU) and baseline (BFS) solutions and evaluating their performance. The next is trying to come up with our own solutions to the problem (PARTIALINDEX, LI, CLUSTEREDEXACT and DOUBLEBFS) and evaluating their performance as well. Finally we demonstrate that our major contribution (LI+OTH+EXTv2) is scalable and can be extended to answer richer queries as well.

There is a public GitHub-repository¹ containing all the code and the thesis.

1.4 Thesis outline

Chapter 2 defines the problem we study in our thesis. Chapter 3 gives a literature overview on “reachability” and “LCR”. Also we discuss some of the ideas for actual algorithms. Chapter 4 explains the actual methods used in the experiments and the concepts behind them. These methods are based on the discussion in the literature review. Chapter 5 discusses the experimental design and the datasets used. Chapter 6 shows the results of the experiments. Chapter 7 concludes on all material and provides work that can be done in the future.

¹<https://github.com/DeLaChance/LCR.git>

Chapter 2

Problem statement

This chapter defines two notions, i.e. “reachability” and “label-constrained reachability” which is abbreviated to “LCR”.

2.1 Reachability in graphs

Definition 2.1.1 defines a graph.

Definition 2.1.1. Graph: A graph $G = (V, E)$ is a pair where V is a set of vertices or nodes and $E \subseteq V \times V$ is a set of edges. The number of vertices is $|V| = n$. The number of edges is $|E| = m$. When a vertex v is in the graph, we say $v \in V$. When an edge from a vertex v to w is in the graph, we say $(v, w) \in E$. In an undirected graph, we have that $(v, w) \in E \Leftrightarrow (w, v) \in E$. For each pair of vertices (v, w) there is at most one edge.

Figure 2.1 shows an example of a graph. The circles are called nodes or vertices and the arrows are called edges. The circles have an number inside them and this is called the vertex identifier (or vertex id).

The notion of “reachability” (Definition 2.2.2) and a “path” (Definition 2.1.2) can formally be defined as:

Definition 2.1.2. Path: Let $G = (V, E)$ be a (directed) graph. Let $s, t \in V$ be two arbitrary nodes. A path P exists in G from s to t if and only if we have one of the following three cases.

1. $s = t$. In this case $P = \langle s \rangle$.
2. $(s, t) \in E$. In this case $P = \langle s, t \rangle$.
3. There exists a sequence of $k \geq 1$ vertices w_1, \dots, w_k such that $(s, w_1) \in E$, for all $1 \leq i \leq k-1$ $(w_i, w_{i+1}) \in E$ and $(w_k, t) \in E$. In this case $P = \langle s, w_1, \dots, w_k, t \rangle$.

Let $\text{Len}(P) \geq 1$ be the length of the path, which is the number of vertices in the path. Let P_i with $1 \leq i \leq \text{Len}(P)$ be the i 'th vertex of the path.

Definition 2.1.3. Reachability: Let $G = (V, E)$ be a (directed). Let $s, t \in V$ be two arbitrary nodes. If there exists a path P from s to t , we say t is reachable from s or $s \rightsquigarrow t$.

The use of reachability queries in graphs has been studied extensively [10, Chapter 6] [1, 9]. A reachability query basically asks whether we can reach from a vertex $v \in V$ another vertex $w \in V$. By this we mean that we can move from v to w in G using any edge in the graph. The vertices that are traversed by going from v to w are said to be a path. In Figure 2.1 vertex 1 can reach vertex 5 over the orange colored vertices. The path that connects 1 to 5 is $\langle 1, 2, 3, 4, 5 \rangle$.

For large scale graphs determining reachability is highly challenging. Certain trade-offs need to be made between index size, query answer time and index construction time. Methods that build up a transitive closure of the graph create an index of size $O(n^2)$ in $O(nm)$ time, but are able to answer a query in $O(1)$ time [10, Chapter 6]. One could think of a transitive closure as a binary matrix of size $n \times n$. Answering these queries on the fly by running a depth- or breadth-first search live which has a running time of $O(n + m)$ but has zero construction time and index size. For large graphs both methods are unacceptable. For example, storing a transitive closure for $n = 1,000,000$ requires 125GB. We need methods that have lower storage requirements at the expense of a higher query time.

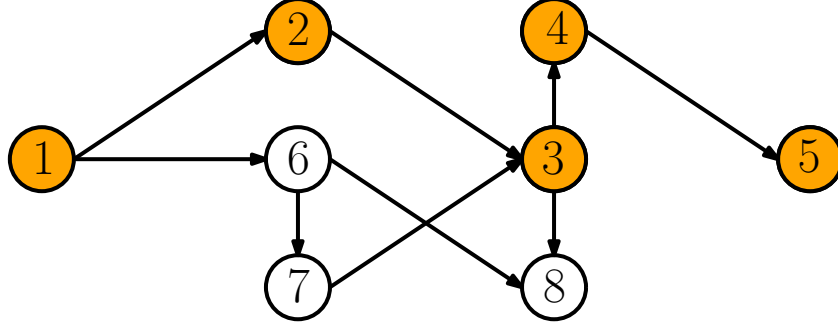


Figure 2.1: An example of a graph with $|V| = 8$ and $|E| = 9$. The orange colored vertices are a path P of the form $\langle 1, 2, 3, 4, 5 \rangle$. The length of the path, i.e. $\text{Len}(P)$, is 5. For each pair of vertices (v, w) with $v, w \in V$ there is at most one edge.

2.2 Label-constrained reachability

First we define a labelled graph. A labelled graph is a graph that has a set of labels, such that each edge $e \in E$ has one specific label.

Definition 2.2.1. Labelled graph: A labelled (directed) graph is a triple $G = (V, E, \mathcal{L})$. \mathcal{L} is a set of labels and $E \subseteq V \times V \times \mathcal{L}$. An edge from $v \in V$ to $w \in V$ can be written as $(v, w) \in E$ or $(v, w, l) \in E$ where $l \in \mathcal{L}$. Moreover let $\text{Label}(e) : E \rightarrow \mathcal{L}$ be a mapping from edge $e \in E$ to its corresponding label. Any $l \in \mathcal{L}$ is called a “label” and any $L \subseteq \mathcal{L}$ is called a “label set”.

In this thesis, we study the problem of label constrained reachability (LCR). Formally this can be defined as:

Definition 2.2.2. Label-constrained reachability (LCR): Let $G = (V, E, \mathcal{L})$ be a (directed) labelled graph. Let $s, t \in V$ be two arbitrary nodes. Let $L \subseteq \mathcal{L}$ be a label set. We say that there is a L -path from s to t if there exists a path P from s to t such that for each edge $(P_i, P_{i+1}, l) \in E$ with $1 \leq i \leq \text{Len}(P)$ we have that $l \in L$. When this is the case, we also write $s \xrightarrow{L} t$.

An example of a directed labelled graph can be seen in Figure 2.2. There exists a $\{a, b\}$ -path from 1 to 5: $\langle 1, 3, 5 \rangle$. From now on when we show a picture of a labelled graph we wish to stick to the following convention: red is label a , blue is label b , green is label c and orange is label d .

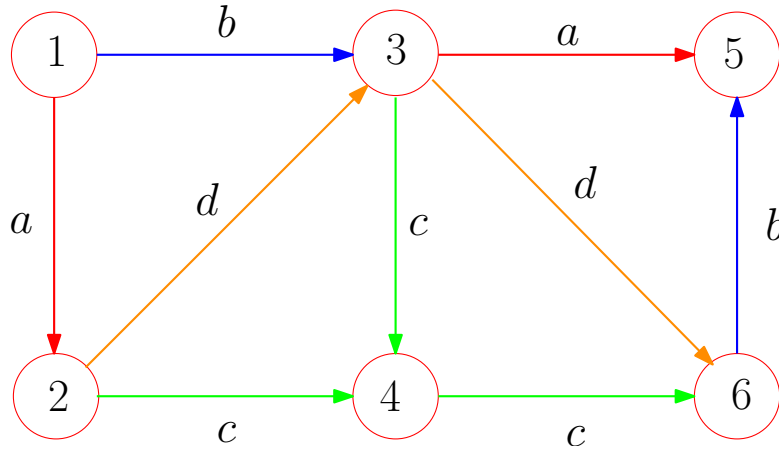


Figure 2.2: An example of a graph with labels. From node 1 to 6 we can see for instance an $\{a, c\}$ -path: $\langle (1, 2, 4, 6) \rangle$. There is also an $\{a, d\}$ -path, a $\{b, c\}$ -path and a $\{b, d\}$ -path. Hence query $q = (1, 6, \{a, c\})$ is answered by true whereas $q = (1, 6, \{c\})$ is answered by false.

LCR was introduced by Jin et al. [5] and further studied by Zou et al. [11] and Chen et al. [2]. The study of LCR was particularly motivated by the study of “regular path queries”. These kinds of queries are prevalent in practical graph query languages such as SPARQL 1.1¹ and Cypher². One can also think of the Facebook-graph (1.35 billion nodes³), where the edge labels indicate a type of relationship (‘friend’, ‘colleague’ or ‘partner’), or Google Knowledge graph (570 million nodes, 18 billion edges⁴), where the edge labels could indicate that two entities are a synonym or entity A is a generalization of entity B, as examples of edge-labelled graphs.

A LCR-query q that determines the existence of a L -path from s to t could be defined as:

Definition 2.2.3. LCR-query: Let $G = (V, E, \mathcal{L})$ be a (directed) labelled graph. Let $s, t \in V$ and $L \subseteq \mathcal{L}$. Then, we can define a query $q \in V \times V \times 2^{\mathcal{L}}$ as (s, t, L) . If $s \stackrel{L}{\sim} t$, then q is said to be true (or a True-query). Otherwise, q is said to be false (or a False-query).

An example query on the graph in Figure 2.2 could be $(1, 6, \{a, c\})$ which would yield true. On the contrary, query $(1, 6, \{c\})$ would be false. The query could also be written as: $(1, 6, (a + c)^*)$, where $(a + c)^*$ indicates that we are allowed to pick any number of a or c ’s in our path from 1 to 6. One could easily think of more enhanced queries like: “Is there a path from 1 to 6 using 4 a ’s first and then at most 3 b ’s?” or “What is the distance in terms of the number of edges or unique labels?”. However, this is not in the scope of our work. We mainly focus on the first type of queries, although we do elaborate on the possibilities to extend our solution to these types of queries.

2.3 The problem

The problem we study is to find a way to answer reachability queries in labelled graphs efficiently and accurately. Trade-offs between index size, query time, index construction time and accuracy have to be made in this regard. Also, we need to take into account that an index might update from time to time due to updates of the graph (node insertion, node removal, edge insertion, edge removal or label change). The frequency of the updates needs to be taken into account.

On the one hand, we can use a slightly adapted version of BFS (or DFS) to answer queries. This approach has the maximal query time but no minimal index construction time. BFS also has minimal update costs.

On the other hand, we can build an index for the full graph that answers all possible queries instantly. This approach has the minimal query time but the maximal index construction time and index size. We assume index construction time and size are correlated. For large graphs building a full index is too cumbersome. Also, updates might require a lot of re-work.

Other approaches will mostly lie into the middle of this spectrum. At the expense of a higher query time, the index time will be reduced.

2.4 Auxiliary definitions

Below we give some auxiliary definitions that will be used along the thesis.

The first definitions (Definitions 2.4.1, 2.4.2, 2.4.3 and 2.4.4) can be used both for labelled and unlabelled graphs. Definition 2.4.1 defines for each vertex $v \in V$ two subsets, the set of vertices that can be reached from v and the set of vertices that can reach v . Definition 2.4.2 defines a strongly connected component (SCC). This is used by ZOU and to define the connectedness of datasets. Definition 2.4.3 defines the property ‘simple’ for a path which is used by Definition 2.4.4. Definition 2.4.4 defines the set of paths between two points, which is used by Definition 2.4.6.

Definition 2.4.1. Ancestors and descendants: In a directed graph for a given node $v \in V$ a set of ancestors and descendants can be defined. Formally, we can define the set of ancestors as $ANCS(v) = \{w \in V \mid w \text{ can reach } v\}$ and the set of descendants $DESC(v) = \{w \in V \mid v \text{ can reach } w\}$. Note that $v \in ANCS(v) \wedge v \in DESC(v)$.

¹<http://www.w3.org/TR/sparql11-query/>, see Section 9: Property paths

²<http://neo4j.com/docs/stable/cypher-query-lang.html>

³<https://nl.wikipedia.org/wiki/Facebook>

⁴<http://www.cnet.com/news/googles-knowledge-graph-tripled-in-size-in-seven-months/>

Definition 2.4.2. Strongly connected component: In a directed graph, a strongly connected component (SCC) is a maximal subset of the nodes $A \subseteq V$ s.t. $\forall [v, v' \in A \mid v \rightsquigarrow v']$. A weakly connected component (WCC) is a subset A s.t. $\forall [v, v' \in A \mid v \rightsquigarrow v' \vee v' \rightsquigarrow v]$.

Definition 2.4.3. Simple path: Let G be a (directed) graph. Let $s, t \in V$. Let P be a path from s to t . P is said to be ‘simple’ if and only if for each P_i with $1 \leq i \leq \text{Len}(P)$ we have that there does not exist P_j with $1 \leq j \leq \text{Len}(P)$ and $j \neq i$ such that $P_i = P_j$.

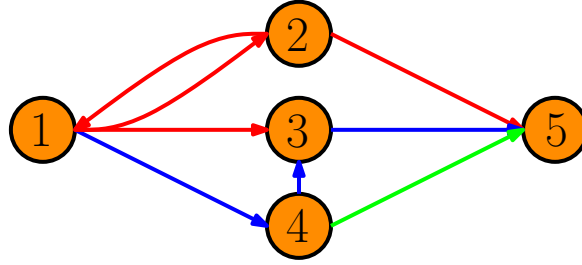
Definition 2.4.4. Set of paths: In a (directed) graph, the set of simple paths between s and t can be expressed as $\text{Paths}(s, t)$. For $s = t$, we have that $\text{Paths}(s, t) = \emptyset$.

Definition 2.4.5 defines the notion of a “path label”, i.e. the union of all labels of edges that are between vertices of the path. Note that there can never be multiple edges between two vertices and hence L_P is unique for any path P . Definition 2.4.6 claims that a path P is minimal if its path label is pairwise incomparable to that of any other path $P' \in \text{Paths}(s, t)$. Definition 2.4.7 defines the set of all minimal paths between s and t . This definition is used by our analysis of algorithm LI. Figure 2.3 shows an example graph G with $\text{Paths}(s, t)$ and $\text{MinPaths}(s, t)$ defined for two particular vertices.

Definition 2.4.5. Path label: In a labelled (directed) graph, let $\text{Labels}(P)$ be defined as: $\bigcup_{i=1}^{\text{Len}(P)-1} \text{Label}(P_i, P_{i+1})$. This is also called “the label of a path” or L_P .

Definition 2.4.6. A minimal path: Let G be a (directed) labelled graph. Let $s, t \in V$. Let $P \in \text{Paths}(s, t)$. P is said to be a minimal path if and only if for all other $P' \in \text{Paths}(s, t)$ we have that $L_P \subseteq L_{P'}$.

Definition 2.4.7. Set of all minimal paths: Let G be a (directed) labelled graph. Let $s, t \in V$. Let $\text{MinPaths}(s, t)$ be the maximal subset of $\text{Paths}(s, t)$ such that for all $P \in \text{MinPaths}(s, t)$ we have that P is minimal.



$$\text{Paths}(1, 5) = \{\langle 1, 2, 5 \rangle, \langle 1, 3, 5 \rangle, \langle 1, 4, 5 \rangle, \langle 1, 4, 3, 5 \rangle\}$$

$$\text{MinPaths}(1, 5) = \{\langle 1, 2, 5 \rangle, \langle 1, 4, 3, 5 \rangle\}$$

Figure 2.3: An example graph G illustrating the definitions. Path $\langle 1, 2, 1, 2, 5 \rangle$ has not been included in $\text{Paths}(1, 5)$ because it is not simple. Path $P_1 = \langle 1, 3, 5 \rangle$ is not included to $\text{MinPaths}(1, 5)$ because of path $P_2 = \langle 1, 4, 3, 5 \rangle$. We have that $L_{P_2} = \{b\}$, whereas $L_{P_1} = \{a, b\}$. Hence path P_1 is not minimal.

Chapter 3

Literature analysis

In this chapter we review a number of papers and chapters that have studied the topics of reachability and LCR. First we review techniques to deal with reachability. Then we go over techniques for answering LCR-queries or some extension of LCR.

3.1 Reachability

In this section we review one method that can answer reachability queries $v \rightsquigarrow w$ for $v, w \in V$ by building an index.

3.1.1 2-hop cover

A 2-hop cover can only be built for a directed acyclic graph (DAG). A DAG can be created out of a directed graph by finding a set of strongly connected components and representing each node as such a SCC (see Definition 2.4.2). Tarjan's algorithm is an algorithm that can do this.

In a 2-hop cover [10, Chapter 6], we assign to each $v \in V$ a 2-hop code which consists of two lists $\text{In}(v) \subseteq \text{ANCS}(v)$ and $\text{Out}(v) \subseteq \text{DESC}(v)$. These lists contain (a subset of) the set of ancestors of v , i.e. $\text{ANCS}(v)$, and the set of descendants of v , i.e. $\text{DESC}(v)$. A valid 2-hop cover is achieved if and only if for every $u, v \in V$ with $u \rightsquigarrow v$ we have that $\text{In}(v) \cap \text{Out}(u) \neq \emptyset$.

A trivial 2-hop cover is the full 2-hop cover. This means that we include for each $v \in V$ all ancestors of v to $\text{In}(v)$ and all descendants of v to $\text{Out}(v)$. This clearly gives no advantage over storing a transitive closure TC . A transitive closure is a $n \times n$ matrix such that $u \rightsquigarrow v$ if and only if entry (i, j) in TC is 1, where i and j are the vertex id's of u and v .

The goal is to find a minimum 2-hop cover. However finding such a cover is NP-hard [10, Chapter 6].



Figure 3.1: A simple graph. In the table below, i.e. Table 3.1, one can see a full hop-cover versus a minimum one.

In Figure 3.1, we can see a simple graph. Its full-hop cover and a minimum cover can be seen in the Table 3.1 below.

Cohen's 2-hop cover approximation

Cohen et al. [3] gives an approximation algorithm which gives a cover at most $O(\log(n))$ larger than the minimum 2-hop cover. The worst case query time is $O(m^{1/2})$ and the index size is $O(nm^{1/2})$.

Algorithm 1 shows this approximation algorithm. Let TC be the transitive closure of G . $(u, v) \in TC$ implies that u can reach v . Consider a node w . Let $A_w \subseteq \text{ANCS}(w)$ and $D_w \subseteq \text{DESC}(w)$. A cluster for w $C_w = S(A_w, w, D_w)$ indicates that every node $u \in A_w$ can reach any node $v \in D_w$. The storage of C_w is

Table 3.1: This table shows in the 2nd and 3rd column a full 2-hop cover and in the 4th and 5th a minimum 2-hop cover.

vertex	full In	full Out	minimum In	minimum Out
1	\emptyset	$\{1, 2, 3\}$	\emptyset	$\{1, 2\}$
2	$\{1\}$	$\{2, 3\}$	$\{1\}$	$\{2, 3\}$
3	$\{2, 3\}$	$\{3\}$	$\{2\}$	\emptyset

$|C_w| = |A_w| + |D_w|$. A cluster C_w indicates that any node in $A_w \subseteq V$ that can reach any node in $D_w \subseteq V$ with a path that goes through w .

Initially we have that for each $w \in V$ $A_w = ASCS(w)$ and $D_w = DESC(w)$ and that $TC' = TC$. As long as TC' is non-empty, i.e. there exists a 1 in the binary matrix of TC' , we continue. We look for a cluster C_w that maximizes the equation on line 5. This is computed by looking at the number of tuples (u, v) where $u \in A_w$ and $v \in D_w$ for which the matrix entry (u, v) in TC is 1. One particular C_w is chosen and for that cluster we remove all (u, v) from TC' and add w to $\text{In}(v)$ and $\text{Out}(u)$.

An interesting observation here is that a full transitive closure is computed at the beginning, whereas one particular reason for using a 2-hop cover was that a transitive closure would use too much memory.

Algorithm 1 2Hop-Cover(G)

```

1: compute  $TC$ 
2:  $TC' \leftarrow TC$ 
3:  $\text{In}(v), \text{Out}(v) \leftarrow \emptyset$  for all  $v \in V$ 
4: while  $TC' \neq \emptyset$  do
5:   find  $\max_{w \in V} |C_w \cap TC'| / (|A_w| + |D_w|)$ 
6:   for  $u \in A_w \wedge v \in D_w$  do
7:     remove  $(u, v)$  from  $TC'$ 
8:     add  $w$  to  $\text{In}(v)$ 
9:     add  $w$  to  $\text{Out}(u)$ 
10:  end for
11: end while
    
```

2-hop cover maintenance

The graph can be updated in four different ways: adding or removing a new node or adding or removing an edge. For labelled graphs also updates of a label may be considered. In [10, Chapter 6], two methods are described to deal with the removal of a node v in a DAG when having a 2-hop code index.

The first method defines a set of nodes $V_{REL} = ANCS(v) \cup DESC(v)$ that consists of all ancestors v and all descendants of v . Let $G_{REL} = (V_{REL}, E_{REL})$ be the subgraph of G where E_{REL} is the maximal subset of E such that for any $(s, t) \in E_{REL}$ we have that $s, t \in V_{REL}$.

A 2-hop cover $L' = (\text{In}', \text{Out}')$ is computed for G_{REL} . Next, we look at all connections $(a, d) \in E$ of which $a \in V_{REL} \vee d \in V_{REL}$. If $a \in V_{REL}$, then we have that $d \in V_{REL}$, $\text{Out}(a) = \text{Out}'(a)$ and $\text{In}(d) = \text{In}'(d)$. If $d \notin V_{REL}$, then $\text{In}(d) = (\text{In}(d) \setminus V_{REL}) \cup \text{In}'(d)$ and $\text{Out}(d) = (\text{Out}(d) \setminus V_{REL}) \cup \text{Out}'(d)$. The downside of this method is the large size of G_{REL} and thereby the high cost of computation.

The second method trades computing time for storage. Suppose we have that $a \in ANCS(v)$ and $d \in DESC(v)$. Let $W \subseteq V$ be the set of all nodes on simple paths from a to d , excluding a and d . Obviously $v \in W$. The trade-off between computing time and storage can be found in storing some nodes $w \in W$ in $\text{Out}(a)$ and $\text{In}(d)$. If node $v \in W$ gets deleted, we can safely delete v in all $\text{Out}(a)$ and $\text{In}(d)$ because there is another route. Of course this method leads to a much lower compression rate.

Applicability to LCR

We review the 2-hop cover method because it can be a part of an algorithm for building an index for LCR-queries. The idea is used by Fletcher and Yoshida [4] in their LCR-index.

LCR adds a whole new dimension to the data. Reachability is like LCR, but then with $|\mathcal{L}| = 1$. Setting $|\mathcal{L}| > 1$ makes it impossible to reduce a group of nodes to a single node like with a SCC without losing information.

A way to use a 2-hop cover for LCR is to first build a 2-hop cover ignoring the edge labels. For any $v \in V$ and $w \in \text{In}(v)$ and for any $v \in V$ and $w \in \text{Out}(v)$ we find $\text{MinPaths}(v, w)$.

3.2 Label-constrained reachability

In this section we discuss two papers [?, 1] and the ideas taken from unpublished notes of Fletcher and Yoshida [4]. The paper by Zou et al. is about building an index for evaluating LCR-queries. The paper by Bonchi et al. [1] is about a more difficult query than a LCR-query. The notes provide us with yet another two algorithms for building an index for LCR-queries.

3.2.1 Bonchi et al.

Bonchi et al. [1] tackle the problem of finding the distance of the shortest path P between two nodes v and w such that $L_P \subseteq L$ for a query (v, w, L) . This is an extension of “LCR”, because a distance is returned if v can reach w and the need to find a shortest distance implies we need to store any path P' with $L_{P'} \supseteq L_P$ if P' is a shorter path than P . Two approaches are given that use landmarks in which the first method is more precise and in which the second is less precise but has a smaller index and needs less construction time.

Both approaches build on the definition of SP-minimality which is similar to Definition 2.4.7. A set of landmarks $X \subseteq V$ is defined. $d_X(s, t)$ denotes the distance between s and t using only edges with labels in X . For simplicity, each edge has a weight of 1. A vertex-pair (u, x) , $u \in V \wedge x \in X$ is said to be SP-minimal w.r.t. a label-set T if and only if there is no $S \subset T$ s.t. $d_T(u, x) = d_S(u, x)$. One should note that the shortest distance can only strictly increase by taking a subset of the labels, because less edges can be used than before.

The first method runs single-shortest path (SSSP) between all landmarks $x \in X \subseteq V$ and all $v \in V$ for a set of candidate label sets $C \subseteq 2^{\mathcal{L}}$. The brute force method takes $O(2^{|C|}k(m+n|C|))$.

However, there are some ways to improve on this running time. Three methods to prune the candidate label sets are mentioned. The first is to look at the labels of edges incident on a landmark x . These labels must be present in a candidate set. The second is the observation that for a given label set L we have that $d_L(x, u) \geq |L|$ because otherwise at least one label $l \in L$ is not used in a path from u to x . The third and last observation uses a history of a set of vertices V_t that are at distance t from the landmark x using a SP-minimal path. The precise details of this last optimization are not entirely clear to us.

The second method assigns a particular label $l \in \mathcal{L}$ to each landmark $x \in X$. Let l_x be that particular label. For each node $v \in V$ we store the “mono-chromatic” distance to each landmark, i.e. $d_{l_x}(v, x)$. The “mono-chromatic” distance is obtained by only using edges of a single type, e.g. only using edges in the label set $\{a\}$. The distances between landmarks are computed using a “bi-chromatic” distance metric, that is using two types of labels e.g. $\{a, b\}$. Given two landmarks $x_1, x_2 \in X$ the distance between x_1 and x_2 would be $d_{\{l_{x_1}, l_{x_2}\}}(x_1, x_2)$. Having k landmarks this can be obtained by doing k breadth-first searches for each node. Hence, we need $O(kn)$ space and $O(k(m+n))$ time. However this solution might return results that are either wrong or far off. For instance, given two landmarks $x_1, x_2 \in X$ we can have that $d_{\{l_{x_1}, l_{x_2}\}}(x_1, x_2) = \infty$ as there is no $\{l_{x_1}, l_{x_2}\}$ -path between the two, but there might be a $\{l_{x_1}, l_{x_2}, l\}$ -path between x_1 and x_2 with $l \in \mathcal{L}$. In this case, a query $q = (x_1, x_2, \{l_{x_1}, l_{x_2}, l\})$ would incorrectly return ∞ .

3.2.2 Zou et al.

In [11] a solution is proposed to solve LCR-queries. It is one of the few paper about “LCR”. The results of the paper are good, in the sense that the index construction time is comparatively low. Hence we have decided to treat it quite detailedly.

In the paper multiple methods for building an index for LCR-queries are being treated. We are only interested in the *transitive-closure method*, as this method has the most promising results. Moreover we only studied the algorithm that builds the index and did not look at the parts about maintenance on which the paper elaborates as well.

We have divided the *transitive-closure method* into 5 steps. First we give an algorithm that can build a transitive closure for a strongly-connected component (SCC). Then, we explain each of the 4 subsequent steps.

From the paper some implementation details remained ambiguous. Moreover there is no proof of correctness for the entire method. Only some parts of the algorithm, e.g. the optimization, have a proof of correctness.

The paper claims that the running time of the full method is $O(\max |V|^3, |V|^2 \cdot D^d)$.

Definitions

The paper has defined some of the definitions in a different way than we do.

Let $G = (V, E, 2^{\mathcal{L}})$ be a directed labelled graph. Each edge $e \in E$ has a label set instead of a label. This is a major difference. Let $\lambda : E \rightarrow 2^{\mathcal{L}}$ be a mapping to a label set.

A ‘single-source transitive closure’ for a graph G is the following: for all nodes $u \in V$ we have that $M_G(u, -)$ can answer any query $q = (u, w, L)$ with $w \in V$ and $L \subseteq \mathcal{L}$ correctly.

Let $P_1, P_2 \in \text{Paths}(s, t)$ for some $s, t \in V$. A path P_2 covers a path P_1 if $L_{P_2} \subseteq L_{P_1}$, where L_P is the ‘path label’ of P (see Definition 2.4.5). The ‘distance of a path P ’ is $|L_P|$.

$\text{Prune}(S)$ is an operator that removes any path P' that is covered by another path P in the set S , essentially creating a minimal path set (see Definition 2.4.7).

\odot concatenates a label set with a set of paths and their associated label sets. For example, given an edge $e = (u, v, \{a\}) \in E$ and a path $P = \langle v, w \rangle$ with $\text{Labels}(P) = \{b\}$, we would have that $\lambda(e) \odot \{P\}$ yields $\{P'\}$, where $P' = \langle u, v, w \rangle$ and $L_{P'} = \{a, b\}$.

Finding single-source transitive closure

First we need an efficient algorithm to find the ‘single-source transitive closure’ for a graph $G = (V, E, 2^{\mathcal{L}})$.

A Dijkstra-like method is proposed to build this index. The algorithm maintains two variables: a heap (or min-priority queue) H and a list RS . H can contain tuples T of the form $(L(P), P, w)$ for $w \in V$. H sorts the tuples according to the number of labels in $L(P)$, i.e. the tuple T for which $|L(P)|$ is minimized is on top of H . We say that a tuple T_1 covers another tuple T_2 if and only if we have that $w_1 = w_2$ and $L(p_1) \subseteq L(p_2)$. RS is initially empty. During the execution of the algorithm it is filled with tuples.

We start at a vertex $u \in V$. Initially, the heap H is filled with all neighbour-tuples of u , i.e. for each $(u, w, L) \in E$ we add an entry of the form $(L, \langle u \rangle, w)$. Then we iterate over the heap taking an entry (or tuple) per iteration. Let $T_1 = (L(P), P, w)$ be the entry taken from the heap. If T_1 is covered by some $T_2 \in RS$ then there is no need to process T_1 and we can continue. Otherwise, we add $L(P)$ to $M_G(u, w)$ and T_1 to RS . After this we can generate neighbour-tuples for w . Any neighbour tuple should have a simple path P and should not already be covered by some tuple on the heap H .

In Figure 3.2 we see a graph for which we could run this algorithm. Let’s say we wish to compute $M_G(u, -)$. First we would generate u ’s neighbour-tuples, which are $(\{b\}, \langle u \rangle, v)$ and $(\{a\}, \langle u \rangle, w)$. Next we would generate for instance the tuple $(\{a, b\}, \{u, w\}, v)$ because it is a neighbour tuple of w . However the tuple $(\{b\}, \{u\}, v) \in RS$ already covers this tuple and hence we do not add this to our index and process the next entry on the heap.

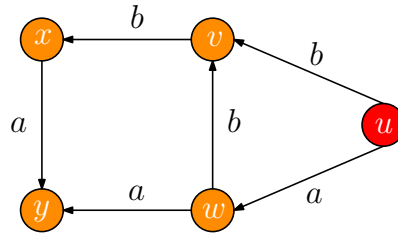


Figure 3.2: A labelled graph as an example. The neighbour-tuples of u are $(\{b\}, \langle u \rangle, v)$ and $(\{a\}, \langle u \rangle, w)$.

The running time of this algorithm is $O(D^d)$ where D is the maximal out-degree and d is the diameter of the graph. This is the theoretical maximum number of paths from any node $u \in V$.

Step 1: Creating a DAG

The first step of Zou's algorithm is to transform a graph G into a DAG D .

First a set k of SCC's $\langle C_1, \dots, C_k \rangle$ is found for the graph G . Each $C_i \subseteq V$ and $\bigcup_{i=1}^k (C_i) = V$. We assume the classical definition of a SCC (see Definition 2.4.2) is used here and not some new definition entailed to labelled graphs. From Figure 3.3 another thing could be assumed, which is merging all vertices $v \in A$ s.t. there is a $\{l\}$ -path between any $v, w \in A$ where $l \in \mathcal{L}$.

Next we generate $M_{C_i}(u, -)$ using the Dijkstra-like method discussed in the previous section for each C_i with $1 \leq i \leq k$.

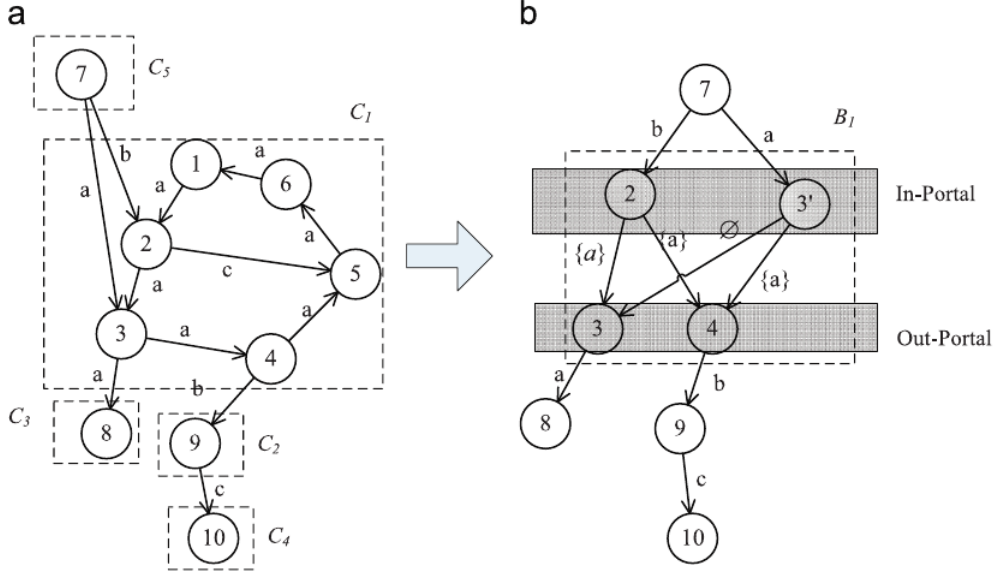


Fig. 4. Augmented DAG. (a) graph G . (b) Augmented DAG D .

Figure 3.3: Figure 4 from paper [11]. It shows a labeled graph G on the left and the 'classical' SCC's on the right.

Next we can generate a DAG $D = (V_D, E_D, 2^{\mathcal{L}})$. Each vertex $v \in C_i$, where C_i is the i 'th SCC, is labelled as either: an in-portal, an out-portal, an internal vertex or both an in- and out-portal. Looking at Figure 3.3, we can say that vertex 6 is an internal vertex and that vertex 3 is both an in- and out-portal. Let $\text{outp}(C_i)$ be the set of out-portals of C_i and let $\text{inp}(C_i)$ be the set of in-portals of C_i .

Let B_i be a bipartite graph. For each $p_1 \in \text{outp}(C_i)$ we add a vertex to B_i . The same is done for each $p_2 \in \text{inp}(C_i)$. p_1 is mapped to one $v \in V$ and p_2 is mapped to some $v' \in V$.

In case $p_1 = p_2$ we add two vertices to B_i . Both map to the same vertex $v \in V$. Also, we add an edge to B_i of the form (p_1, p_2, \emptyset) .

In case $p_1 \neq p_2$ we look for all label sets L in $M_{C_i}(p_1, p_2)$. For each such label set we add an edge of the form (p_1, p_2, L) . Figure 3.4 shows an example of a transformation of a SCC C_i to a bipartite graph B_i .

By merging all bipartite graphs B_i we get D . In Figure 3.5 you can see an example of the full process. Note that the edges in D are label sets and not labels.

One can be curious as to how effective this approach would be in graphs in which (almost) every node in the graph is an in- and/or out-portal. Then there needs to be an edge from every node in the SCC to every other node, which is quadratic in the number of nodes in the SCC. There might be even more edges in D than in G . Thus, the effectiveness of this approach really depends on the percentage of nodes that are serving as either in- and/or out-portals.

Figure 3.4 illustrates this case. On the left a part of a graph eligible to become a SCC is displayed. On the right the result. We see that the number of edges has increased a lot.

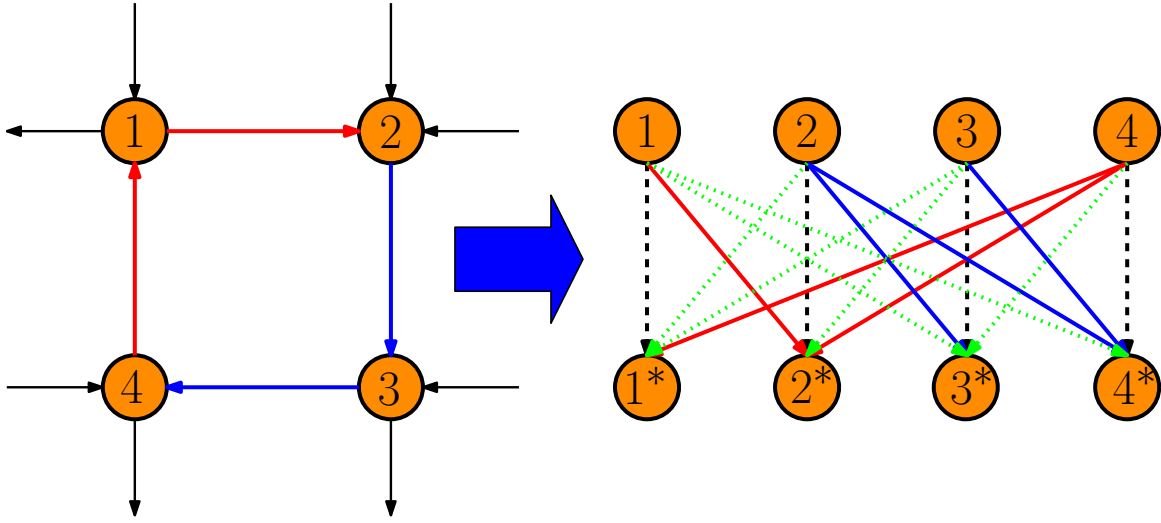


Figure 3.4: A part of G eligible to become a SCC C_i on the left. Nodes 1 to 4 form a SCC. On the right we see the resulting bipartite graph B_i . A dashed line indicates an edge with no label set, because the end points of that edge represent the same vertex in C_i . A colored edge indicates a connection between two portals in C_i that are connected through edges with only that color. A dotted green edge indicates the portals are connected through a combination of red and blue edges. We see that the new number of edges is much higher than in the original.

Step 2: using D

As D is a DAG we can generate a reversed topological order RT for it. The internal vertices of any SCC C_i are not in D . Only the vertices that serve as an in- or out-portal of any C_i are in D .

Let $u \in RT$. First we set for each node $M_G(u, -)$ to $M_{C_i}(u, -)$ if $u \in C_i$. For each out-edge $(u, v, L) \in E_D$ we do the following. We set $M_G(u, -)$ to $M_G(u, -) \cup \text{Prune}(\lambda(u, v) \odot M_G(v, -))$. This takes all the entries (w, L') from $M_G(v, -)$ for $w \in V$ and tries to insert $(w, L' \cup L)$ into $M_G(u, -)$.

Figure 3.5 illustrates the idea. The figure displays the graph as it is, but you might as well see the set of bipartite graphs D as displayed on the right part of the figure. Changes are propagated from bottom to top in this graph.

We start at vertex $12 \in C_3$ after having computed its internal transitive closure. As 12 has no children, we go to C_2 . For each out-portal in C_2 (which is $11'$) connected to an in-portal in C_3 (which is 12), we concatenate all entries of the local transitive closure of the in-portal with the edge label of the edge between that in- and out-portal. In the case of C_2 this would set $M_{C_2}(11', 12)$ to $\text{Prune}(\{b\} \odot M_{C_3}(12, -) = \{b\}$. For each in-portal of C_2 (which are vertices 10 and 11) we look at all the other out-portals in C_2 (which is $11'$). In the case of C_2 and vertex 10 this would set $M_{C_2}(10, 11')$ to $\text{Prune}(\{a\} \odot M_{C_2}(11', -)) \cup \text{Prune}(\{b\} \odot M_{C_2}(11', -))$, as there is an $\{a\}$ and a $\{b\}$ -path connecting 10 and 11 in C_2 . We do not need to do anything for vertex 11, as $M_{C_2}(11, -) = M_{C_2}(11', -)$. After we have completed this process for C_2 we do the same for C_1 .

Step 3: Inner to out

Now we have only covered the vertices $v \in V_D \subseteq V$, i.e. we have not covered the internal (or inner) vertices of each SCC. To get $M_G(u, -)$ for all nodes $u \in V$ we need to look to the in- and out-portals of each SCC respectively.

The inner vertices are processed in an order PT that adheres to RT . In case the in- and out-portals of C_j were processed before the in- and out-portals of C_i in RT , then the inner vertices of C_j will be processed before the inner vertices of C_i in PT .

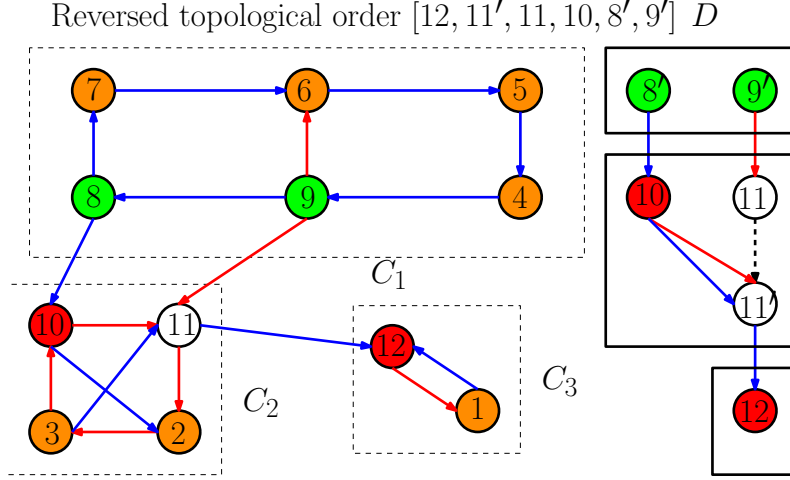


Figure 3.5: On the left a graph with SCC's C_1 , C_2 and C_3 , which could also be represented by their related bipartite graphs. The edge labels are indicated by the color. A dashed line indicates an edge with no label, which happens when two nodes in D represent the same vertex in G . Inner vertices are orange, out-portals are green and in-portals are red and nodes serving as both are white. The internal vertices are numbered. On the right the graph D holding all bipartite components. Changes are propagated upwards along all edges in the graph where the transitive closure of a node is concatenated with the edge label.

For each inner vertex $v \in PT$ and out-portal $p \in C_i$ we set $M_G(v, -)$ to $\bigcup_{p \in \text{outp}(C_i)} (\text{Prune}(\lambda(v, p) \odot M_G(p, -)))$.

Step 4: Final step

The final step is then to iterate over all C_i . For any $v \notin C_i$, $u \in C_i$ and $u' \in \text{outp}(C_i)$ we set $M_G(v, u) = \text{Prune}(M_G(v, u') \odot M_G(u', u)) \cup M_G(v, u)$. In this way any inner vertex of C_j 'knows about' the inner vertices in C_i .

We did not understand the necessity of this final step. In step 2 the in- and out-portals $p \in C_i$ learn about all inner vertices $v \in C_j$ in case the in- and out-portals of C_j preceded those of C_i in RT . In step 3 the inner vertices $v \in C_i$ take the entries of the out-portals $p \in \text{outp}(C_i)$. If we have that for $v \in C_i$ and $w \in C_j$ $v \rightsquigarrow w$, then v should 'know about' w at the end of step 3.

Optimization

The paper has an optimization to this algorithm which improves the results dramatically. For instance, for a 100k graph the optimization could reduce the index construction time from 5,000 seconds to about 15 seconds. It tries to optimize the computation of $M_{C_i}(u, -)$.

A minor change to the Dijkstra-like algorithm is made. Suppose we wish to compute $M_{C_i}(u, -)$ for $u \in V$. Moreover we have already computed $M_{C_i}(v, -)$ for another $v \in V$ and $(u, v) \in E$, then we can start by setting $M_{C_i}(u, -)$ to $M_{C_i}(u, -) \cup \{\lambda(u, v) \odot M_{C_i}(v, -)\}$.

Zou results

The code was written in C++ and the experiments were conducted on a P4@3.0Ghz machine 2 GB RAM running Linux Ubuntu.

A number of synthetic datasets were generated using either the Erdos-Renyi (ER) or the Scale-Free (SF) model. ER is a classical random graph. $|E|$ edges are chosen randomly from $|V|(|V| - 1)$ possible edges. Each edge in this model is equally likely which may make it not comparable to real-life networks. SF graphs try to include the notion of preferential attachment, i.e. a node with a lot of connections (in and/or out) may be more likely to be connected to or from. Hence, we get a skewed out-degree distribution of the graph. The out-degree

distribution of such a graph follows the shape of an exponential distribution and can be specified by the following 4 parameters: $|V|$, the minimal degree of a node, the maximal degree of a node and γ^1 . Typically, $2 \leq \gamma \leq 3$.

We have that $|\mathcal{L}| = 18$. The number of labels in any query was $|\mathcal{L}| \cdot 0.3$. This means that each query (s, t, L) has $|L| = 6$.

In Figure 3.6, we can see the results for ER-graphs of different sizes. $|d| = 1.5$ means that the average degree per node was 1.5. We see that the optimized version of the *transitive-closure method* performs much better. Looking at the last column we can see that the query times obtained by this method are roughly equal to the query times obtained by doing a double-sided BFS, i.e. a BFS with two threads.

Performance VS. $ V $ in ER Graphs.					
1.5	Transitive closure method				Bi-directional search
	IT (s)	IT-opt (s)	IS (KB)	QT (ms)	QT (ms)
	15	3	415	0.01	0.01
	23	5	1396	0.01	0.02
	217	7	4920	0.02	0.02
	623	10	9000	0.03	0.03
	964	12	13,200	0.03	0.04
2	5065	15	33,000	0.04	0.05

Figure 3.6: Table 4 from paper [11]. For different sizes of ER-graphs it shows the time to build the graph using the *transitive-closure method* and the optimized version of it, the size in KB and the query time in ms. In the last column we see the query times obtained by a double-threaded BFS.

In Figure 3.7 we see the effects of increasing the density in ER-graphs. The number of vertices has been set to $|V| = 10,000$.

In Figure 3.8 we see the results for SF-graphs. There is a major difference here compared to 3.6. The paper argues that most nodes in these graphs have a very low degree and hence much of the search space can be pruned

¹<http://tinyurl.com/zzt8qgq>

Table 5
Performance VS. Density in ER Graphs.

Degree	Transitive closure method				Bi-directional search
d	IT (s)	IT-opt (s)	IS (KB)	QT (ms)	QT (ms)
2	6890	20	26.3	0.08.	0.05
3	11,112	23	102.3	0.09	0.09
4	25,347	35	160	0.12	0.15
5	33,169	80	186	0.23	0.18

Figure 3.7: Table 5 from paper [11]. For different densities of 10k ER-graphs it shows the time to build the graph using the *transitive-closure method* and the optimized version of it, the size in KB and the query time in ms. Interestingly the index size is very small. In the last column we see the query times obtained by a double-threaded BFS.

Table 7
Performance VS. $|V|$ in SF graphs.

$ V $	Transitive closure method				Sampling-tree method			Bi-directional search
$d=1.5$	IT (s)	IT-opt (s)	IS (KB)	QT (ms)	IT (s)	IS (KB)	QT (ms)	QT (ms)
1K	0.1	0.1	43	0.01.	0.6	1576	0.16	0.01
2K	0.2	0.1	94	0.01	2.6	2583	0.17	0.02
4K	0.2	0.1	140	0.01	9.5	4870	0.18	0.02
6K	0.4	0.2	289	0.01	27.8	8854	0.20	0.03
8K	0.4	0.2	390	0.01	45.5	11393	0.22	0.03
10K	1.1	0.3	492	0.01	80.4	23931	0.24	0.03

Figure 3.8: Table 7 from paper [11]. For different sizes of SF-graphs it shows the time to build the graph using the *transitive-closure method* and the optimized version of it, the size in KB and the query time in ms.

early on.

One curious thing to note here is that the index size of any synthetic graph in these experiments is relatively low. Most do not exceed a megabyte. In our experiments (see Table 6.3) almost any index would have a size that exceeds this. In Figure 3.6 we see that an ER-graph with 10k vertices and a degree of 1.5 (i.e. 15k edges) has an index with a size of 33,000 KB. The 10k ER-graphs in Figure 3.7 with a higher degree all have an index with a lower size, e.g. for $d = 5$ we see an index with a size of 186KB.

Looking at Figure 3.7 we see a 300 to 400 speed-up in the index construction time between the normal (IT) and the optimized version (IT-opt). We have implemented a similar optimization for our major contribution (see Section 4.2.2). However we have not observed an improvement of this magnitude for the index construction time.

3.2.3 Fletcher and Yoshida

There are two possible algorithms taken from unpublished notes of Fletcher and Yoshida [4] for determining reachability in labelled graphs. One algorithm DOUBLEBFS performs two breadth-first searches (BFS'es) from each node $v \in V$. The other algorithm is an iteration-based algorithm in which nodes send messages to their neighbours.

DoubleBFS and NeighbourExchange

In this section we describe two methods (DOUBLEBFS and NEIGHBOUREXCHANGE) that can be used for building a full exact index.

Both algorithms try to create two arrays $\text{Out}(v)$, $\text{In}(v)$ similar to the idea of a 2-hop cover. We define $\text{In}(v) = \{(s, L) \mid \exists [P \in \text{MinPaths}(v, s) \mid \text{Labels}(P) = L]\}$ and $\text{Out}(v) = \{(s, L) \mid \exists [P \in \text{MinPaths}(s, v) \mid \text{Labels}(P) = L]\}$. The formula to answer a reachability query $q = (s, t, L)$ in the labelled-graph is:

$$Q(s, t, L) = \exists [v \in V \mid \exists [L_s, L_t \subseteq \mathcal{L} \mid (v, L_s) \in \text{Out}(s) \wedge (v, L_t) \in \text{In}(t) \wedge L_s \subseteq L \wedge L_t \subseteq L]].$$

First, some ordering of the vertices $\langle v_1, \dots, v_n \rangle$ is created. This ordering has little or no effect on the algorithms discussed here, but it might have in more sophisticated versions of these algorithms later on.

The first algorithm loops over all n vertices in the ordering that has been created. From each v_i , we start two Breadth-First searches (BFS'es): one using the direction of the edges and one using the opposite direction. Each time we hit a node (v, L) we try to add it to $\text{Out}(v_i)$ or $\text{In}(v_i)$ respectively.

The second algorithm NEIGHBOUREXCHANGE starts with a loop that continues as long as there is a node v which had a change in its $\text{In}(v)$ or $\text{Out}(v)$. During each such iteration we loop over all v_i with $1 \leq i \leq n$. A vertex $v_i \in V$ can send an update (x, L) to any $w \in V$, where $x \in V$ is the source of the update and $L \subseteq \mathcal{L}$ is the propagated label set. Each v_i first processes its updates. For any applied update (x, L) , i.e. it was inserted into the index of v_i , and for any neighbour (v_i, w, l) with $l \in \mathcal{L}$ v_i sends an update to w of the form $(x, L \cup \{l\})$.

Both algorithms use a form of pruning. Suppose we wish to add (u, L) to either $\text{In}(v_i)$ or $\text{Out}(v_i)$ and we have that $Q(v_i, u, L)$ is already true, we do the following. DOUBLEBFS would not push any of u 's neighbours on the queue. NEIGHBOUREXCHANGE would not push the update (u, L) to any of v_i 's neighbours.

In Figure 3.9 and Table 3.2, we can see how algorithm NEIGHBOUREXCHANGE works. In the table we see the messages that are being exchanged during the first round and how that effects the $\text{In}(v)$ of each node $v \in V$.

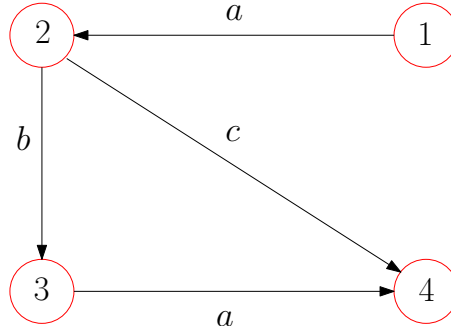


Figure 3.9: A graph with labels. In Table 3.2 we see the updates during the first round of the algorithm NEIGHBOUREXCHANGE.

While working on the code for DOUBLEBFS and NEIGHBOUREXCHANGE and running experiments for it, we came along some issues.

Firstly, NEIGHBOUREXCHANGE and DOUBLEBFS build both In and Out , but having one of these suffices to answer any query directly. This halves the memory usage. Secondly, NEIGHBOUREXCHANGE and DOUBLEBFS do a lot of redundant computation. Suppose we have a graph like in Figure 3.10 in which the black-box in the middle could be any kind of graph which makes the red nodes reachable from the blue nodes. For each blue node we have to redo a BFS through the black box on the red nodes. We might as well compute $\text{Out}(A)$ first and then re-use it each time a BFS hits A .

Table 3.2: This table shows the first and only round of running NEIGHBOUREXCHANGE on the graph in Figure 3.9. A node first processes its updates and then sends new updates to its out-neighbours for each change that was made.

Vertex-id	Receive	Send
1	\emptyset	$(1, \{a\})$
2	$(1, \{a\})$	$(1, \{a, b\}), (2, \{b\}), (2, \{a, c\}), (2, \{c\})$
3	$(1, \{a, b\}), (2, \{b\})$	$(1, \{a, b\}), (2, \{a, b\})$
4	$(1, \{a, c\}), (2, \{c\}), (2, \{a, b\}), (3, \{a\})$	\emptyset

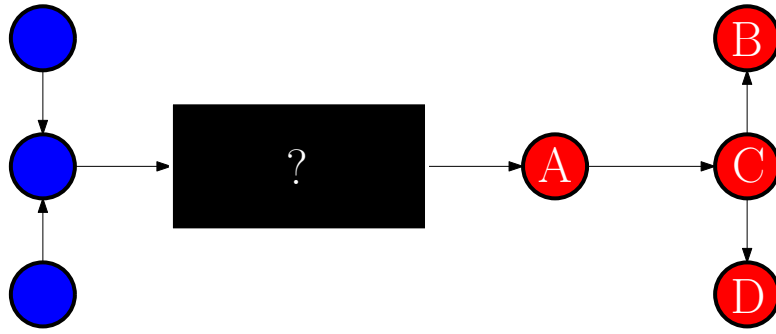


Figure 3.10: A graph. The black box could be any graph which makes the red nodes reachable from the blue nodes.

Managing the index size: union or intersection

The index size could become a real bottleneck for larger $|V|$ when building a full exact index. A solution to this by Fletcher and Yoshida [4] is to not add all entries to the index, but instead merge the label sets of some of the entries.

Suppose we have an entry of the form $(w, \{\{a, b\}\}, (w, \{b, c\}), (w, \{c, d, e\}))$ in $\text{In}(v)$ for some $v, w \in V$. We could try to reduce the size of the index by intersecting the label sets or taking the union of the label sets. A union would yield $(w, \{a, b, c, d, e\})$. An intersection would yield $(w, \{b\})$. A query $(v, w, \{b, d\})$ would return false in the first case and true in the second case, whereas we know the query is a False-query considering the original entry. On the other hand, the True-query $(v, w, \{b, c\})$ would return false in the first case and true in the second case. Hence, a union might produce false negatives, whereas an intersection might produce false positives.

Merging the label sets should only be done whenever the number of label sets of an entry exceeds some budget $b \geq 0$. As soon as we add a new label set L_2 to an entry containing b label sets, we merge it with an existing label set L_1 . In case we use a union-policy, we merge L_2 with a label set L_1 such that $|L_1 \cup L_2|$ is minimal. In case we use an intersection-policy, we merge L_2 with a label set L_1 such that $|L_1 \cap L_2|$ is minimal. Moreover, we give the merged label sets a flag such that we know it is a merged entry.

Chapter 4

Methods

In this chapter we discuss the methods either used in at least one experiment or that we have at least strongly considered using (e.g. JOINDEX). Also we discuss index maintenance and extensions to other queries which are richer than LCR-queries. For our major contribution (LI) we give a proof of correctness and a worst-case running time and memory usage analysis.

4.1 Existing methods

This section discusses existing methods to answer LCR-queries.

4.1.1 BFS

BFS is the baseline idea. It has no index. The time to construct is the ‘index’ is equal the time to load the graph into memory. The same holds for the size of the ‘index’. Updates are as simple as adding or removing a node or vertex from the graph. In many cases BFS is a good solution, as it is relatively fast, requires no construction and maintenance and has a low memory usage.

Algorithm 2 shows the precise implementation for a query (v, w, L) . Just like normal BFS, this BFS has a running time complexity of $O(|V| + |E|)$.

Algorithm 2 BFS(v, w, L)

```
1: Let marked :  $V \rightarrow \{0, 1\}$  be a mapping with for all  $v \in V$  marked[ $v$ ] = 0
2: Let  $q$  be a queue
3:  $q.push(s)$ 
4: while  $q$  is not empty do
5:    $v \leftarrow q.pop()$ 
6:   marked[ $v$ ]  $\leftarrow 1$ 
7:   if  $v = w$  then
8:     return True
9:   end if
10:  for  $(v, w, l) \in E \wedge \text{marked}[w] = 0$  do
11:    if  $l \in L$  then
12:       $q.push(w)$ 
13:    end if
14:  end for
15: end while
16: return False
```

The variable `marked` at line 1 has been implemented using a bitset, which makes it possible to modify and check for any value in constant time rather than in $O(\log(n))$ at the expense of $O(n)$ memory usage.

4.1.2 Best effort Zou

This method is a best effort approach to implement Zou’s algorithm [11].

Let G be a labelled directed graph (Definition 2.2.1). Let $M_G(u, -) \subseteq V \times 2^{\mathcal{L}}$. This is the same as $\text{Ind}(u)$ in our implementation. Hence, one might as well replace $M_G(u, -)$ by $\text{Ind}(u)$. The operation $M_G(u, -) = \text{Prune}(M_G(u, v) \odot M_G(v, -))$ is implemented by looking at any $(v, L) \in \text{Ind}(u)$ and concatenating (v, L) to any $(w, L') \in \text{Ind}(v)$ which results in $(w, L \cup L')$. $(w, L \cup L')$ is attempted to be added to $\text{Ind}(u)$ by using TRYINSERT (Algorithm 11, Section 4.2.2).

The algorithm consists of a number of steps:

1. Generate a set of k SCC’s $C^* = \{C_1, \dots, C_k\}$, by using Tarjan’s algorithm for G .
2. For each $C_i \in C^*$ we build an index $M_{C_i} = \bigcup_{u \in C_i} (M_{C_i}(u, -))$ answering all queries local to C_i . After this step, for any $u, v \in C_i$ with a L -path from u to v we have that there exists an entry (v, L') in $M_{C_i}(u, -)$ with $L' \subseteq L$.
3. We find the in- and out-portals of each SCC C_i , i.e. we find those vertices $v \in C_i$ such that there exists an incoming edge $(w, v) \in E$ or an outgoing edge $(v, w) \in E$ where $w \notin C_i$. We define these sets as $\text{IN}(C_i)$ and $\text{OUT}(C_i)$ respectively.
4. We generate an acyclic graph $D = (V_D, E_D, \mathcal{L})$.

The vertices $v \in V_D$ are in- and out-portals of a SCC C_i in G . In case a vertex $v \in C_i$ is both an in- and out-portal we generate a replica vertex v' which is added to V_D . Between any such v and v' there is an edge: $(v, v', \{\emptyset\})$.

The edges $e \in E_D \subseteq V \times V \times 2^{\mathcal{L}}$ have label sets rather than labels onto them. Let $\lambda : E \rightarrow 2^{\mathcal{L}}$ be a function for D that maps each edge to a particular label set. Between any two $v, w \in V_D$ there can be multiple edges, but these need to have different label sets.

For each $v \in \text{IN}(C_i)$ and $w \in \text{IN}(C_i)$ with $v \neq w$ we look for all label sets L in $M_{C_i}(v, w)$ and generate an edge of the form (v, w, L) . This edge is added to D .

5. D is acyclic and can be topologically sorted. This is done, after which we reverse the order to create an order RT .
6. For $u \in RT$ and all children $v \in V_D$ of u we set $M_G(u, -)$ to $\text{Prune}(\bigcup_{(u,v) \in E_D} (\lambda(u, v) \odot M_{C_i}(v, -)))$. This puts an entry $(w, \lambda(u, v) \cup L_2)$ in $M_G(u, -)$, if there is an entry (w, L_2) in $M_G(v, -)$.
7. For each out-portal u_i that is part of C_i we look for all inner vertices u_j of any C_j with $j \neq i$ and set $M_G(u_i, -) = \text{Prune}(M_G(u_i, u_j) \odot M_G(u_j, -)) \cup M_G(u_i, -)$. In this way the out-portal ‘knows about’ the inner vertices of other SCC’s as well. This step was not in the original paper, but necessary to get our implementation fully working in the end.
8. For each inner vertex u_i in C_i we look at all out-portals p_i of C_i and set $M_G(u_i, -)$ to the union of all out-portals, i.e. we set $M_G(u_i, -)$ to $\bigcup_{p_i \in \text{OUT}(C_i)} (\text{Prune}(M_G(u_i, p_i) \odot M_G(p_i, -)))$. This puts any entry $(u_j, L_1 \cup L_2)$ in $M_G(u_i, -)$, if there is an entry (p_i, L_1) in $M_G(u_i, p_i)$ and an entry (u_j, L_2) in $M_G(p_i, u_j)$.
9. For each inner vertex u_j with $j \neq i$ we look for an in-portal p_i such that $u_j \rightsquigarrow p_i$ and any other vertex $u_i \neq p_i$. We set $M_G(u_j, -)$ to $\text{Prune}(M_G(u_j, p_i) \odot (M_G(p_i, u_i)) \cup M_G(u_j, -)$.

We think steps 7 and 9 are redundant, even though step 9 is in the paper. Given two vertices $u \in C_i$ and $v \in C_j$ such that $u \xrightarrow{L} v$ we should have that there exists some out-portal $u' \in \text{OUT}(C_i)$ such that $L \in M_G(u', v)$ at the end of step 6. Step 8 should then include L to $M_G(u, v)$ via u' . This did not happen in some rare cases. Hence we included step 7 to get a fully working version.

We decided not to try to fix these issues or implement the optimization described in the paper, because the results were already quite disappointing, i.e. they were not competitive to e.g. DOUBLEBFS and not in line with the results in the original paper. The difference in index construction time between DOUBLEBFS and ZOU was at least a factor 10. We did not see how this gap could be bridged. Steps 2, 7, 8 and 9 took the most time. The approach works very bad for graphs for which there is at least one relatively large SCC. Chapter 6 elaborates on ZOU with experimental results as well.

We included the idea of step 2 of the algorithm into our implementations of DOUBLEBFS and LI, i.e. sort the entries of the heap according to the number of labels in the label set.

4.2 Our contributions

In this section we discuss our contributions for answering LCR-queries. The approach LI with its extensions is our major contribution.

4.2.1 General comments

Let $\text{Ind} \subseteq V \times V \times 2^{\mathcal{L}}$ be an index for all nodes $v \in V$. Let $\text{Ind}(v) \subseteq V \times 2^{\mathcal{L}}$ be the index for node v and let $\text{Ind}(v, w) \subseteq 2^{\mathcal{L}}$ be a list of (minimal) label sets connecting v and w .

$\text{Ind}(v)$ is said to be minimal if for any two different pairs (u, L) and (u, L') in $\text{Ind}(v)$ we have that neither $L \subseteq L'$ or $L' \subseteq L$. An index Ind is said to be minimal if and only if for all $v \in V$ we have that $\text{Ind}(v)$ is minimal.

4.2.2 LandmarkedIndex

While writing and testing the code we found that BFS had a low query answering time, e.g. on a $|V| = 1,000$ -graph the query answering time would often be in an order of magnitude of 10^{-6} . We also saw that building the full index takes a lot of time and memory. Building a full index is often not needed to achieve a speed-up. For instance, if a sufficiently large subset of the nodes $V' \subseteq V$ is (fully) indexed, we already can achieve a speed-up. Hence, we decided to opt for an approach that is similar to BFS but uses an index.

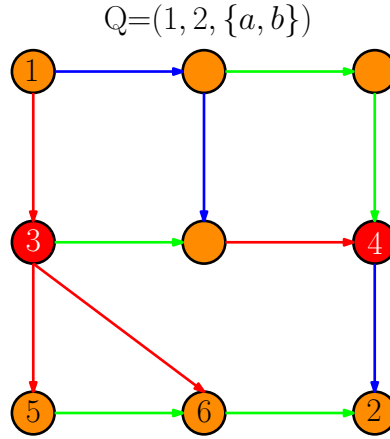


Figure 4.1: The query is $(1, 2, \{a, b\})$. The red edges are a , the blue edges b and the green edges c . Also let the red vertices 3 and 4 be landmarks. The query is a True-query, as 1 can reach 4 and 4 can reach 2.

LandmarkedIndex builds a complete index for a subset $V' \subseteq V$ of $|V'| = k$ nodes, also called *landmarks*. All the remaining vertices are called *non-landmarks*. At query evaluation time, we run a normal BFS for a query (v, w, L) . However, when we hit a node v' that is a landmark, we try to resolve the query directly. If the answer is true, we are done. If the answer is false, we do not need to process any of the out-neighbours of v' saving time compared to BFS.

We introduce three extensions to the basic LI approach. Table 4.1 shows the name and the present extensions of each version of LI that was used in at least one experiment. Table 5.2 in Chapter 4 gives a summary of all methods.

Basic idea

First we select k landmarks. This can be done in various ways, but the default way is picking the k vertices with the highest total degree. k has been specified before and is typically a fraction of n . After this we create a mapping $\text{isL} : V \rightarrow \langle -1, 0, \dots, k-1 \rangle$ that maps any vertex $v \in V$ to either -1 if it is not a landmark and a positive integer in the range $0 \dots k-1$ if it is a landmark. We say $v \in V'$ or $v_i \in V$ for $i \geq 1$ if v is a landmark and $v \in V \setminus V'$ or $v_0 \in V$ if v is not a landmark.

Table 4.1: This table shows the name and the present extensions of each version of LI that was used in at least one experiment.

Method name	Ext. 1	Ext. 2	Ext. 3
LI			
LI+EXTv1	✓		
LI+OTH		✓	
LI+OTH+EXTv1	✓	✓	
LI+OTH+EXTv2		✓	✓

The basic approach, i.e. without any extensions, runs Algorithm 3 for any landmark node first. The $i + 1$ 'th call of the algorithm pushes a pair $(v_i, \{\})$ to the queue initially. Then, it iterates until the queue is empty. For each entry (u, L) obtained on line 4 we call TRYINSERT. It tries to insert (u, L) into $\text{Ind}(v_i)$ and returns true if and only if this succeeded. TRYINSERT only succeeds whenever L is incomparable to any other L' in $\text{Ind}(v_i, u)$ or whenever L is only a subset of some of the entries in L' . Moreover, it removes any superset L' of L in $\text{Ind}(v_i, u)$. Hence, TRYINSERT preserves the minimality of $\text{Ind}(v_i, u)$. Finally, we add all pairs $(w, L \cup \{l\})$ to q for each $(u, w, l) \in E$. The purpose of FORWARDPROP is explained in the following section.

Algorithm 3 LabelledBFSPerNode(v)

```

1: Let  $q$  be a queue
2:  $q.\text{push}(v, \{\})$ 
3: while  $q$  is not empty do
4:    $(u, L) \leftarrow q.\text{pop}()$ 
5:   if  $\text{tryInsert}(u, v, L) = \text{False}$  then
6:     continue
7:   end if
8:   if  $\text{hasBeenIndexed}(u) = 1$  then
9:      $\text{forwardProp}(v, u, L)$ 
10:    continue
11:  end if
12:  for  $(u, w, l) \in E$  do
13:     $q.\text{push}(w, L \cup \{l\})$ 
14:  end for
15: end while
16:  $\text{hasBeenIndexed}(v) \leftarrow 1$ 

```

Algorithm 4 $\text{tryInsert}(u, v, L)$

```

1: if  $v = u$  then
2:   return True
3: end if
4: if  $\exists [(u, L') \in \text{Ind}(v) \mid L' \subseteq L]$  then
5:   return False
6: end if
7: remove any  $(u, L')$  s.t.  $L \subseteq L'$  from  $\text{Ind}(v)$ 
8: add  $(u, L)$  to  $\text{Ind}(v)$ 
9: return True

```

Algorithm 5 forwardProp(v, u, L)

```

1: for  $(w, L') \in \text{Ind}(u)$  do
2:   tryInsert( $v, w, L \cup L'$ )
3: end for

```

Propagation

In our code we used by default ‘forward propagation’. While building the index for a node v_j with $j \geq i$, we might visit a node v_i which has already been full indexed. When there is a L_1 -path from v_j to v_i , we can say that any entry $(u, L_1 \cup L_2)$ could be appended to $\text{Ind}(v_j)$ where $(u, L_2) \in \text{Ind}(v_i)$. The idea came from studying the optimization of Zou’s algorithm [11]. Forward propagation is used by the basic approach LI.

During construction we maintain a mapping that indicates which vertices $v \in V$ have already been fully indexed by using $\text{hasBeenIndexed} : V \rightarrow \{0, 1\}$. $\text{hasBeenIndexed}(v_i) = 0$ if and only if landmark v_i has not been fully indexed. Otherwise, we have that $\text{hasBeenIndexed}(v_i) = 1$.

Heap

Algorithm 3 can explore parts of the graph multiple times. This can hurt the performance. Inspired by Zou et al. [11] we changed the queue into a heap (or a min priority queue) where the heap entries are sorted on the number of labels in the label set. The entry (u, L) for which $|L|$ is minimal is on top of q . In this way we are assured that we never insert (u, L) before (u, L') to $\text{Ind}(v_i)$ when $L \supset L'$. Hence TRYINSERT never has to remove any entries. The heap might incur some overhead, but on the other hand might save time when a large part of the graph is traversed multiple times. The heap is used by the basic approach LI.

Basic query algorithm

Algorithm 6 is the pseudocode for resolving a query for LI+OTH. To get the version without any of the extensions you should ignore lines **5-13**.

We start by exploring the graph using a BFS at line **14**. Only when we hit a landmark $v' \in V'$ during this BFS, we try to resolve the query directly. This is done by calling Algorithm 7. This looks for an entry of the form $(w, L') \in \text{Ind}(v')$. It returns true if and only if such an entry exists and $L' \subseteq L$. When this ‘direct attempt’ fails we do not need to look at any of the out-edges of v' , because we know w cannot be reached from v' .

First and third extension

In Figure 4.1 we see that vertex 1 can reach landmark 3, but 3 cannot reach vertex 2 over $\{a, b\}$. However, 3 can reach two other vertices: vertices 5 and 6. If one of these vertices were able to reach 2 then 3 could reach 2. Hence there is no point in evaluating any vertex that has a $\{a, b\}$ -path from vertex 3.

This notion can be generalized. The basic idea for the first extension is the following. Given a query $q = (v, w, L)$ we might stumble upon a landmark v' one way or another. If a direct attempt (v', w, L) fails, we know that any vertex w' that has a L -path from v' cannot reach w . Hence it can be pruned. We also call a query that prunes in this way an ‘extensive query’.

When we use the first extension, we choose to store some extra information for each landmark $v' \in V$. We store $l' = |\mathcal{L}|$ extra lists containing all vertices that can be reached using one particular label only. If a query (v', w, L) fails from a landmark v' , we can prune all vertices that are in one of the in total l' lists if the corresponding label of that list is in L .

Let $\text{LP}(v) : \mathcal{L} \rightarrow 2^V$ be the variable holding these lists for each $v \in V'$. LP is filled during each iteration of Algorithm 3 after line **7**. We obtain an entry (u, L) at line **4**. As long as $L = \{l\}$ and $u \neq v$ we insert u to $\text{LP}(v)(l)$.

One could think of maintaining more lists, i.e. also adding lists representing the all label sets with two or three labels. However this can be quite expensive in cost of memory as there can be $\binom{|\mathcal{L}|}{3}$ such lists, which is $O(|\mathcal{L}|^3)$. There is also going to be a lot of redundancy in these lists.

As the first extension did not provide enough speed-up for False-queries for when $k \leq \frac{n}{50}$ (see Chapter 6,

Algorithm 6 queryLM(v, w, L)

```

1: if isL( $v$ )  $\geq 0$  then
2:   return queryDirect( $v, w, L$ )
3: else
4:   Let  $\text{marked} : V \rightarrow \{0, 1\}$  be a mapping
5:   for  $v' \in \{v' \mid (v', L) \in \text{Ind}(v)\}$  do
6:     //  $v'$  is always a landmark
7:     for  $L' \in \text{Ind}(v, v') \wedge L' \subseteq L$  do
8:       if queryDirect( $v', w, L$ ) = True then
9:         return True
10:      end if
11:       $\text{marked}[v'] \leftarrow 1$ 
12:    end for
13:  end for
14:  Let  $q$  be a queue
15:   $q.\text{push}(v)$ 
16:  while  $q$  is not empty do
17:     $u \leftarrow q.\text{pop}()$ 
18:     $\text{marked}[u] \leftarrow 1$ 
19:    if  $u = w$  then
20:      return True
21:    end if
22:    if isL( $u$ )  $\geq 0$  then
23:      if queryDirect( $u, w, L$ ) = True then
24:        return True
25:      end if
26:    continue
27:    end if
28:    for  $(u, u', l) \in E \wedge \text{marked}[u'] = 0$  do
29:      if  $l \in L$  then
30:         $q.\text{push}(u')$ 
31:      end if
32:    end for
33:  end while
34:  return False
35: end if

```

Algorithm 7 queryDirect(v, w, L)

```

1: //  $v$  is a landmark
2: for  $L' \in \text{Ind}(v, w)$  do
3:   if  $L' \subseteq L$  then
4:     return True
5:   end if
6: end for
7: return False

```

Section 6.4) we needed to revisit the idea of the first extension. This extension can only be used as a substitute for the first extension and it uses the same principle. It is either the first or the third extension, but never both.

For each landmark $v' \in V$ we maintain $l' \leq 2^{D_{\max}}$ entries with $D_{\max} = |\mathcal{L}|/4 + 1$. The variable holding these entries for a landmark $v' \in V'$ is called $\text{seqE}(v') \subseteq 2^V \times 2^{\mathcal{L}}$ in our code. Each entry $(L, B) \in \text{seqE}$ consists of a set $B \subseteq V$ and a label set L . Each B is a ‘bitset’ of n bits in which a bit j is set if and only if we can reach the corresponding vertex w which has vertex id j . Hence reading or writing for a specific bit takes $O(1)$ time.

The third extension inserts some extra code to Algorithm 3 after line 8. Each time an entry (u, L) has been

discovered for $\text{Ind}(v)$ on line 4, we can do two things. Either, we create a new entry $(L, B) \in V \times \mathcal{L}$ where $B = \{u\}$. Or, we find the already existing entry (L, B) and set B to $B \cup \{u\}$.

The third extension also inserts some extra code to Algorithm 3 after line 15. Given two entries $(L_1, B_1), (L_2, B_2) \in \text{seqE}$ where $L_1 \subseteq L_2$ we set B_2 to $B_2 \cup B_1$. After this is done for all entries for all entries in $\text{seqE}(v)$. We sort the entries $(L, B) \in \text{seqE}$ descendingly w.r.t. $|L|$.

The advantage of the third extension over the first is that we can use larger label sets, i.e. from 1 to D_{\max} . We chose to create a maximal distance, as we found that medium-size label sets can often cover a large subset of the nodes in graphs with a high degree or an exponential label set distribution, and can be used a majority of the queries we generated. The second advantage is that we can quickly join the bitset `marked` used by the query algorithm (Algorithm 6) with the bitset B used by any entry $(L, B) \in \text{seqE}(v)$. Joining can be done in $\frac{n}{64}$ time on a 64-bit machine.

Second extension

After all landmark vertices have been fully indexed, the second extension can add a number of entries (v_i, L) (where v_i is a landmark) to $\text{Ind}(v_0)$ for all non-landmark vertices $v_0 \in V$.

We choose b to be a constant, typically a low value. Each v_0 adds at most b entries to $\text{Ind}(v_0)$. After it has found this many paths or it cannot find any more entries, it stops.

The algorithm for doing this is very similar to Algorithm 3. The differences are listed below.

1. We only add entries (v_i, L) to $\text{Ind}(v_0)$ with $i \geq 1$.
2. We stop immediately after having successfully inserted b entries to $\text{Ind}(v_0)$.
3. We use both landmarked and non-landmarked nodes for FORWARDPROP, i.e. we copy the entries from each indexed node $v \in V$ regardless of whether it is a landmark. However in case $v \in V'$ we only look at entries $(v_j, L) \in \text{Ind}(v)$ with $j \geq 1$. `hasBeenIndexed` is set to 1 at the end for any $v_0 \in V$ as well.
4. We add a variable `marked` : $V \rightarrow \{0, 1\}$ before line 3. Each entry (u, L) found on line 4 sets `marked(u)` $\leftarrow 1$. In case we have that `marked(u)` = 1 we continue, i.e. we skip processing entry (u, L) .

Query algorithm with extensions

Adding the second extension to LI enables lines 5-13 of Algorithm 6. In these lines we look for all entries $(v_i, L') \in \text{Ind}(v)$ with $i \geq 1$.

The difference between LI+OTH and LI+OTH+EXTv1 is that the first call of QUERYDIRECT has been replaced by a call to QUERYEXTENSIVEDIRECT (in case v is not a landmark). The difference between LI+OTH and LI+OTH+EXTv2 is that any call QUERYDIRECT has been replaced by a call to QUERYEXTENSIVEDIRECT.

Algorithm 8 consists of two versions. One is intended for the first extension and one is intended for the third extension. Both prune all w' that can be reached from a landmark v' with label set L that could not reach the query-target w with that same label set.

Space complexity

We analyze the worst-case memory usage for LI+OTH+EXTv1 and LI+OTH+EXTv2. k is the number of landmarks and b is the budget per non-landmark node.

Considering only the landmarks, the memory usage can in the worst case be $O(|V| \cdot k \cdot 2^{|\mathcal{L}|})$. Each landmark $v' \in V'$ can index up to $|V| - 1$ other nodes and for each such node we can store at most $2^{|\mathcal{L}|}$ paths.

Each landmark v' (k in total) can also have $|\mathcal{L}|$ lists in LP and each such list can be of length n at most. This results in $O(|\mathcal{L}| \cdot n)$ extra storage for v' .

Each landmark v' can also have $2^{D_{\max}} = 2^{|\mathcal{L}|/4+1}$ entries in `seqE`. Each consists of a set of at most $O(n)$ (B) and a constant $O(1)$ (L). This results in $O(2^{|\mathcal{L}|} \cdot n)$ extra storage for $v' \in V$.

Each remaining non-landmark vertex $v_0 \in V$ ($n - k$ in total) can store up to b entries, which takes $O(b)$.

Hence in total the memory usage is: $O(n \cdot (b + k2^{|\mathcal{L}|}))$ (LI+OTH+EXTv1 and LI+OTH+EXTv2).

Algorithm 8 queryExtensiveDirect(v, w, L, marked)

```

1: // first extension
2: if queryDirect( $v, w, L$ ) = True then
3:   return True
4: end if
5: for  $l \in L$  do
6:   for  $v' \in \text{LP}(v)(l)$  do
7:      $\text{marked}(v') \leftarrow 1$ 
8:   end for
9: end for
1: // third extension
2: if queryDirect( $v, w, L$ ) = True then
3:   return True
4: end if
5: for  $(L', B) \in \text{seqE}$  do
6:   if  $L' \subseteq L$  then
7:      $\text{marked} \leftarrow \text{marked} \cup B$ 
8:     break
9:   end if
10: end for

```

Index construction time complexity

We analyze the worst-case running time for LI+OTH+EXTv1 and for LI+OTH+EXTv2.

For each non-landmark $v_0 \in V$ we can visit all vertices $v \in V$ and make either b (non-landmark) or k (landmark) TRYINSERT-calls.

For each landmark $v' \in V'$ we can have that TRYINSERT returns true at most $2^{|\mathcal{L}|}$ times. Hence each edge can be visited by each landmark at most that many times. As we also visit all the vertices building the index for each landmark takes at most $O(k(n+m)2^{|\mathcal{L}|})$.

For a landmark $v' \in V'$ filling $\text{seqE}(v')$ takes at most $2^{|\mathcal{L}|}$ time per call to TRYINSERT, as in the worst case we need to loop over at most $2^{|\mathcal{L}|}$ entries to find a specific entry. Adding a vertex $u \in V$ to a set V^* for an entry $(L', V^*) \in \text{seqE}$ can be done in constant time.

For a landmark $v' \in V'$ inserting a certain $u \in V$ for a pair $(u, \{l\})$ to $\text{LP}(v')(l)$ takes $O(\log(n))$ per call to TRYINSERT.

The index construction time complexity is $O(n^2 2^{|\mathcal{L}|} \cdot (b+k) + kD^d 2^{|\mathcal{L}|})$ (LI+OTH+EXTv2) and $O(n^2 2^{|\mathcal{L}|} \cdot (b+k) + kD^d \log(n))$ (LI+OTH+EXTv1).

Query answering time complexity

The running time of Algorithm 7, that is QUERYDIRECT, is $O(2^{|\mathcal{L}|} + \log(n))$ as there are at most $2^{|\mathcal{L}|}$ label sets between any two vertices which need to be compared. Finding any specific w in $\text{Ind}(v)$ takes at most $\log(n)$ time.

The running time of Algorithm 8, that is QUERYEXTENSIVEDIRECT, is that of QUERYDIRECT plus the worst case time of doing at most $2^{|\mathcal{L}|}$ comparisons and setting at most n bits in marked . Hence this takes $O(n + 2^{|\mathcal{L}|})$.

In case of LI+OTH+EXTv1 or LI+OTH+EXTv2, the total worst case running time of Algorithm 6 is $O(n + m + k \cdot (2^{|\mathcal{L}|} + n))$. At most k direct attempts to resolve a query are made, either on lines **5-13** or on lines **23-25**. The graph exploration part of Algorithm 6 takes at most $O(n + m)$.

Overview of running time and memory usage analysis

Table 4.2 shows an overview of the worst-case index construction time, index size and query answering time for LI+OTH+EXTv1 and LI+OTH+EXTv2.

Table 4.2: This table shows an overview of the worst-case index construction time, index size and query answering time for LI+OTH+EXTv1 and LI+OTH+EXTv2.

Description	LI+OTH+EXTv1	LI+OTH+EXTv2
Worst-case index size	$O(n \cdot (b + k2^{ \mathcal{L} }))$	$O(n \cdot (b + k2^{ \mathcal{L} }))$
Worst-case index construction time	$O(n^2 2^{ \mathcal{L} } \cdot (b + k) + k(n + m)2^{ \mathcal{L} })$	$O(n^2 2^{ \mathcal{L} } \cdot (b + k) + k(n + m)2^{ \mathcal{L} })$
Worst-case query answering time	$O(n + m + k \cdot (2^{ \mathcal{L} } + n))$	$O(n + m + k \cdot (2^{ \mathcal{L} } + n))$

Proof of correctness

Let $v, w \in V$ and let $L \subseteq \mathcal{L}$. We say $query(v, w, L) = \text{True}$ if and only if there is a L -path from v to w . Otherwise we say $query(v, w, L) = \text{False}$. We begin by proving that $query(v, w, L) = \text{True} \Leftrightarrow$ there exists $(w, L') \in \text{Ind}(v)$ with $L' \subseteq L$.

Proposition 1. *Let $V' \subset V$ be the set of landmark vertices. Let $\langle v_1, \dots, v_k \rangle$ be an ordering of the k landmark vertices. Let $\text{Ind}(v_j)$ be the index of the j 'th landmark constructed by Algorithm 3. We wish to establish that for any $w \in V$ (1) $query(v_j, w, L) = \text{True}$ implies that there exists $(w, L') \in \text{Ind}(v_j)$ with $L' \subseteq L$ and (2) $query(v_j, w, L) = \text{False}$ implies that there does not exist $(w, L') \in \text{Ind}(v_j)$ with $L' \subseteq L$. We do this by induction on $1 \leq j \leq k$.*

Proof. **Base case:** $j = 1$.

Algorithm 3 has only been called for v_1 in this case.

(1): Assume $query(v_j, w, L) = \text{True}$. Then there is a L -path P from v_j to w . Lines **12-16** have found any label set (minimal or not) connecting v to some $u \in \text{DESCS}(v_j)$, where $\text{DESCS}(v_j)$ is the set of descendants of v_j . Only line **5** could have blocked adding (w, L) to $\text{Ind}(v_j)$ in this case. Let (u, L') be the first vertex in path P where the if-condition on line **5** is true, i.e. TRYINSERT returns false. This only happens if there exists $(u, L^*) \in \text{Ind}(v_j)$ with $L^* \subseteq L'$. This argument can be repeated for any next vertex in path P and hence we have that there exists $(w, L') \in \text{Ind}(v_j)$ where $L' \subseteq L$.

(2): Assume $query(v_j, w, L) = \text{False}$. Then there is no L -path P from v_j to w . If there is no such path lines **12-16** will never push an entry of the form (w, L') s.t. $L' \subseteq L$. Hence such an entry can never appear in $\text{Ind}(v)$.

Step: $1 < j \leq k$. We assume that for v_1, \dots, v_{j-1} we have that **(1)** and **(2)** are true (IH).

Algorithm 3 has been called for $\langle v_1, \dots, v_j \rangle$ in this case.

(1): Assume $query(v_j, w, L) = \text{True}$. Then there is a L -path P from v_j to w . Lines **12-16** have found any label set (minimal or not) connecting v to some $u \in \text{DESC}(v_j)$ (see Definition 2.4.1), where $\text{DESC}(v_j)$ is the set of descendants of v . W.r.t. line **5** we have the same argument as in the base case. Line **8** could have been true for any v_i with $i \leq j - 1$. Let v_i (with label set $L_1 \subseteq L$) be the first vertex on path P such that this happens. In that case by the induction hypothesis we have that there exists $(w, L_2) \in \text{Ind}(v_i)$ where $L_2 \subseteq L$. The entry $(w, L_1 \cup L_2)$ will be added to $\text{Ind}(v_j)$.

(2): Assume $query(v_j, w, L) = \text{False}$. Then there is no L -path P from v_j to w . If there is no such path lines **12-16** will never push an entry of the form (w, L') s.t. $L' \subseteq L$. As **(2)** holds for any v_i with $i \leq j - 1$ an incorrect entry is not added by any of the v_i to $\text{Ind}(v_j)$. \square

Corollary 4.2.0.1. *We wish to establish that $query(v_j, w, L) = \text{True} \Leftrightarrow$ there exists $(w, L') \in \text{Ind}(v_j)$ where $L' \subseteq L$.*

Proof. We have that either $query(v_j, w, L) = \text{True}$ or $query(v_j, w, L) = \text{False}$. From Theorem 1 we get that $query(v_j, w, L) = \text{False} \Rightarrow$ there does not exist $(w, L') \in \text{Ind}(v_j)$ where $L' \subseteq L$. By contraposition this implies that there exists $(w, L') \in \text{Ind}(v_j)$ where $L' \subseteq L \Rightarrow query(v_j, w, L) = \text{True}$. \square

Corollary 4.2.0.2. *After building the index for all non-landmark vertices $v_0 \in V$, we have that for landmark $v_i \in V$ with $i \geq 1$ there exists $(v_i, L') \in \text{Ind}(v_0)$ with $L' \subseteq L \Rightarrow query(v_0, v_i, L) = \text{True}$.*

Proof. Let $v_0 \in V$ be an arbitrary non-landmark node. The entries $(v', L') \in \text{Ind}(v_0)$ where v' is a landmark obtained for non-landmark nodes $v_0 \in V$ are created by an algorithm with the same structure as Algorithm 3. There are four differences compared to Algorithm 3: at most b entries can be added to $\text{Ind}(v_0)$, only entries of the form (v_i, L') with $i \geq 1$ are added, each vertex u in an entry (u, L) is visited at most once on line 4 of Algorithm 3 and both landmark and non-landmark vertices are used for forward pruning.

The differences are restrictions on the number and the type of those entries that can be added to $\text{Ind}(v_0)$. The entries that satisfy the restrictions, are still valid entries. Hence, by Proposition 1 and Corollary 4.2.0.2 we get that there exists $(v_i, L') \in \text{Ind}(v_0)$ with $L' \subseteq L \Rightarrow query(v_0, v_i, L) = \text{True}$. \square

Next we wish to verify the correctness of QUERYDIRECT and QUERYEXTENSIVEDIRECT.

Proposition 2. *QUERYDIRECT returns true if and only if $query(v, w, L) = \text{True}$. QUERYEXTENSIVEDIRECT only prunes a subset of the vertices $V^* \subseteq V$ s.t. for all $v^* \in V^*$ we have that $query(v^*, w, L) = \text{False}$.*

Proof. By Corollary 4.2.0.1 we know that QUERYDIRECT return true if and only if $query(v, w, L) = \text{True}$.

Variables LP and seqE used by the first and third extension respectively are filled by using entries (u, L') obtained in the loop in lines 3-16 of Algorithm 3. As each such entry is in $\text{Ind}(v)$ and $query(v, w, L) = \text{True} \Leftrightarrow$ there exists $(w, L') \in \text{Ind}(v)$ with $L' \subseteq L$ by Corollary 4.2.0.1, we have that any vertex u included in seqE and LP is reachable from v using L' .

Lines 5-9 (first extension) or 5-10 (third extension) are only reached in case QUERYDIRECT returns false on line 2. Any vertex $v^* \in V^*$ set marked on line 7 (both extensions) can be reached by some subset $L' \subseteq L$ from v . If such v^* were able to reach w over L then v could reach w as well. Hence such v^* cannot reach w and can correctly be pruned, i.e. included to marked. \square

Finally we wish to prove the correctness of QUERYLM.

Theorem 4.2.1. *Let $v, w \in V$ and $L \subseteq \mathcal{L}$. We have that algorithm 6 (using either LI, LI+OTH, LI+OTH+EXTv1 or LI+OTH+EXTv2) returns true if and only if $query(v, w, L) = \text{True}$.*

Proof. This can be proven by a case distinction.

- **Case 1:** On line 2 we run QUERYDIRECT for (v, w, L) . Because of Proposition 1 and Proposition 2, the statement holds.
- **Case 2:** On line 8 we run QUERYDIRECT (or QUERYEXTENSIVEDIRECT) for (v', w, L) where v' is a landmark that can be reached from v using L (Corollary 4.2.0.2). Because Proposition 2 and Proposition 1, the statement holds and we prune a subset of the correct vertices, i.e. those vertices that have no L -path to w , in case of QUERYEXTENSIVEDIRECT.
- **Case 3:** On line 8 we run QUERYDIRECT (or QUERYEXTENSIVEDIRECT) for (v', w, L) where v' is a landmark that can be reached from v using L . The v to v' path has been discovered in the loop from lines 16-33. Because of Proposition 1, the statement holds. Because of Proposition 2 we prune a subset of the correct vertices in case of QUERYEXTENSIVEDIRECT.
- **Case 4:** In this case the loop on lines 16-33 ended. There cannot be a L -path from v to w as no QUERYDIRECT (or QUERYEXTENSIVEDIRECT) call found such a path, only correct vertices were pruned according to Proposition 2 and this path was not found through graph exploration (Lines 16-33). The algorithm returns false.

\square

DoubleBFS

DOUBLEBFS is a special case of LI, i.e. $k = n$ where the number of landmarks k equals the number of vertices n . Any extension is irrelevant in this case as any query is answered by QUERYDIRECT.

4.2.3 Partial

PARTIALINDEX also builds an index for a subset of k nodes like LI, but does not build a full index to cut down the index size. Rather it chooses a budget b which is a fraction of n , e.g. $\frac{n}{4}$. Each of the k nodes then index up to b nodes fully ignoring any other vertices. The difference with LI is that at query evaluation time we have to evaluate the out edges of a node v' . Otherwise, we might return wrong results.

PARTIALINDEX has a frequency. The frequency determines the number of nodes that needs to be hit before we try to resolve a query directly. Having to resolve a lot of direct queries can slow down the query evaluation and also often have no purpose if the sources of the direct queries are close to each other. Figure 4.2 shows an example. In the final implementation the frequency is roughly $O(\sqrt{n})$.

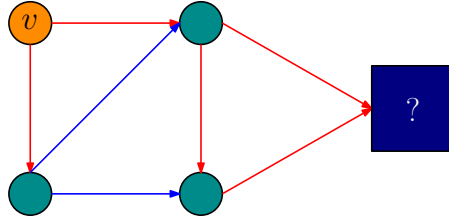


Figure 4.2: Suppose we start a query at orange vertex v . Let the box on the right be a large graph and the cyan-colored vertices be vertices for which there is a partial index. Each of these also has a partial index. Having to run a direct query for each of these will most likely yield the same result.

4.2.4 NeighbourExchange

This method is based on the discussion in Chapter 2 Section 3.2.3 with some differences applied. The algorithm builds a full index. Algorithm 9 shows a pseudocode.

On lines **3** to **7** we initialize **Upd** with the direct neighbour updates. This is necessary to kick off the entire algorithm. Next, we start a loop at line **8** that continues as long as there are updates to the index. We loop over any node $v \in V$ and check whether there are entries in **Upd**[v]. We add an entry (w, L) on line **16** to **Ind**(v) if and only if it does not violate the minimality of the index. Next, we add the neighbours of w .

One of the disadvantages of NEIGHBOUREXCHANGE is that it has to use **Upd** and **Ind**. The memory usage of **Upd** can grow considerably depending on the order in which the updates are processed. Another disadvantage is that we see no way to use propagation for this index.

4.2.5 Joindex

This method is based on the discussion in Section 3.2.3 with some differences applied. JOINDEX can both be used with intersection or union.

The algorithm to construct JOINDEX is very similar to that of DOUBLEBFS or NEIGHBOUREXCHANGE. We need two (exact) indices **ln** and **Out**. First, we need to determine a budget b which is the maximum number of entries per pair (u, v) for some $u, v \in V$ and decide whether we wish to use intersection or union to merge the entries. For union we do the following. While adding an entry (u, L) to **Ind**(v), we verify whether this would result in exceeding b . If this is the case, we compare L to all b entries L' that are in the index. The entry (u, v, L') for which $|L' \cup L|$ is minimized, is merged with (u, v, L) yielding an entry $(u, v, L \cup L')$ for which a join flag is set

Algorithm 9 NeighbourExchange()

```

1:  $c \leftarrow \text{True}$ 
2: let  $\text{Upd} \subseteq V \times V \times 2^{\mathcal{L}}$ 
3: for  $v \in V$  do
4:   for  $(v, w, l) \in E$  do
5:     add  $(v, \{l\})$  to  $\text{Upd}[w]$ 
6:   end for
7: end for
8: while  $c = \text{True}$  do
9:    $c \leftarrow \text{False}$ 
10:  for  $v \in V$  do
11:     $d \leftarrow \text{False}$ 
12:    for  $(w, L) \in \text{Upd}[v]$  do
13:      if  $\text{tryInsert}(v, w, L) = \text{True}$  then
14:         $c \leftarrow \text{True}$ 
15:         $d \leftarrow \text{True}$ 
16:      end if
17:      if  $d = \text{True}$  then
18:        for  $(v, w', l) \in E$  do
19:          if  $l \in L$  then
20:            add  $(v, \{l\} \cup L)$  to  $\text{Upd}[w']$ 
21:          end if
22:        end for
23:      end if
24:    end for
25:     $\text{Upd}[v] \leftarrow \emptyset$ 
26:  end for
27: end while

```

indicating it has been merged (the most significant bit). For intersection, we do the same except the resulting entry is $(u, v, L \cap L')$. At query evaluation time, the join flag can be used.

Algorithm 10 shows the query algorithm in case JOININDEX uses union to merge entries. A query (v, w, L) is resolved by looking for a node s.t. $u \in \text{Out}(v) \cap \text{In}(w)$. If we have that $(u, L) \in \text{Out}(v)$ and $(u, L) \in \text{In}(w)$ and neither has a join flag on, we return True. Only in case the left or right part (u, L') has a join flag, indicated by $J(u, L')$, and the entry is a superset of (u, L) , we need to recur on that part. If we have tried out all u , then we return False. The algorithm for answering JOININDEX queries using intersection is very similar.

In the end we decided not to fully implement JOININDEX, for the following reasons.

1. JOININDEX needs two times the index, In and Out, and hence has a double memory consumption for at least some time period.
2. JOININDEX only decreases the index size but not the construction time. LI and PARTIALINDEX decrease both the construction time and index size.
3. JOININDEX increases the query evaluation time. NEIGHBOUREXCHANGE or DOUBLEBFS can answer a query in $O(\log |V| + 2^{|\mathcal{L}|})$, whereas JOININDEX needs $O(|V| + |E|)$.

4.2.6 ClusteredExact

This method is based on the observation that separating the graph into some subgraphs and building a full index for each of these subgraphs individually is more efficient than for the whole graph.

First, we create a clustering. Each of the N nodes is assigned to one of the K clusters. An exact index is generated, e.g. by using DOUBLEBFS, for each cluster.

Algorithm 10 JoindexQueryUnion(v, w, L, marked)

```

1: if  $|L| = 0$  then
2:   return False
3: end if
4: if  $v = w$  then
5:   return True
6: end if
7: for  $u \in (\text{Out}(v) \cap \text{In}(w) \wedge \text{marked}[u] = 0)$  do
8:   if  $\exists [(u, L') \in \text{Out}(v) \mid L' \subseteq L]$  then
9:     if  $\exists [(u, L') \in \text{In}(w) \mid L' \subseteq L]$  then
10:      return True
11:    else
12:      if  $\exists [(u, L') \in \text{In}(w) \mid L' \supseteq L \wedge J(u, L')]$  then
13:         $\text{marked}[u] \leftarrow 1$ 
14:        JoindexQueryUnion( $u, w, L, \text{marked}$ )
15:         $\text{marked}[u] \leftarrow 0$ 
16:      end if
17:    end if
18:  else
19:    if  $\exists [(u, L') \in \text{Out}(v) \mid L' \supseteq L \wedge J(u, L')]$  then
20:       $\text{marked}[u] \leftarrow 1$ 
21:      JoindexQueryUnion( $v, u, L, \text{marked}$ )
22:       $\text{marked}[u] \leftarrow 0$ 
23:    end if
24:  end if
25: end for
26: return False

```

The query evaluation algorithm is the most tricky part about this approach. Algorithm 11 gives a pseudocode. $\text{clD} : V \rightarrow \langle 0, \dots, K-1 \rangle$ maps each vertex to a certain cluster. outP gives the list of vertices that are out ports for that cluster. Line 6 checks whether the query can be resolved directly because x and w are in the same cluster. Otherwise, we loop over all out-ports p of the cluster x belongs to. It is checked whether x can reach p using L and whether the label l of an edge (p, s) is in L . If this is the case, we push s onto the stack.

Algorithm 11 ClusteredExactQuery(v, w, L)

```

1: Let marked :  $V \rightarrow \{0, 1\}$ 
2: Let  $q$  be a queue
3: while  $q$  is not empty do
4:    $x \leftarrow q.pop()$ 
5:   if  $clD[x] = clD[w]$  then
6:     if  $directQuery(x, w, L) = \text{True}$  then
7:       return True
8:     end if
9:   end if
10:  if  $marked[x] = 1$  then
11:    continue
12:  end if
13:   $marked[x] \leftarrow 1$ 
14:  for  $p \in outP[clD[x]]$  do
15:    if  $directQuery(x, p) = \text{True}$  then
16:      for  $(p, s, l) \in E \wedge l \in L$  do
17:        push  $s$  onto  $q$ 
18:      end for
19:    end if
20:  end for
21: end while

```

4.3 Implementation details

Let $G = (V, E, \mathcal{L})$ be an arbitrary labelled directed graph.

There is a difference between labels and label sets. A label is an individual element in \mathcal{L} , e.g. $a \in \{a, b, c\}$. Labels are represented by a number in the range $0 \dots |\mathcal{L}| - 1$. Label sets are represented as bitsets which means that if $L \subseteq L'$ we mean that in our code L and L' have a non-empty bit-intersection, e.g. $L = \{a, b\} = 3 = (00000011)$ and $L' = \{a, b, c\} = 7 = (00000111)$ have a non-empty bit-intersection namely the first and second bit whereas $4 = (00000100) = \{c\}$ and $2 = (00000010) = \{b\}$ have none. 0 denotes the empty label set.

The code has the ability to scale down the number of bits used in a label set, e.g. if $|\mathcal{L}| = 8$ we can use just a byte, thereby cutting down the index size. However for all experiments we set the size of a label set to 4 bytes. We did this because there were some datasets in our experiments with $|\mathcal{L}| \geq 8$ and a few with $|\mathcal{L}| \geq 16$.

An index **Ind** can have two modes: blocked mode and non-blocked mode. In blocked mode we store the index as a triple-array with dimensions: $|V| \times |V| \times |2^{\mathcal{L}}|$, in which the third dimension can scale depending on the actual number of label sets for a given pair (v, w) . In non-blocked mode we store the index for each node $v \in V$ as a list of tuples of the form: (w, L^*) where L^* is a list of label sets. Blocked mode can cut down the index construction time compared to non-blocked mode, but takes much more memory and is not scalable. Hence we only used non-blocked mode in our experiments.

4.4 Index maintenance

In this section we stress the issue of index maintenance. Suppose that after running LI, PARTIALINDEX, DOUBLEBFS, CLUSTEREDEXACTOR NEIGHBOUREXCHANGE we have built an index **Ind**(v) for every $v \in V$. There can be 5 update events:

1. Addition of a node v
2. Removal of a node v
3. Addition of an edge (v, w, l)
4. Removal of an edge (v, w, l)

5. Changing edge label (v, w, l) to (v, w, l')

A choice we made for all approaches is to build a minimal index, i.e. for each pair of nodes (v, w) we have that all label sets L and L' are no super- or subset of each other. This can have implications for index maintenance, as the removal of one edge can necessitate rebuilding a part of the index. This part can become quite large. Like discussed in Section 3.1.1 storing some redundant information might help making maintenance easier.

Figure 4.3 illustrates this problem. Suppose we would remove edge $(2, 3, \{a\})$, then we need to update 0 and add entries $\{a, b\}$ and $\{a, c\}$ to $(0, 4)$ while removing $\{a\}$. For graphs with a lot of cycles this problem can get even worse.

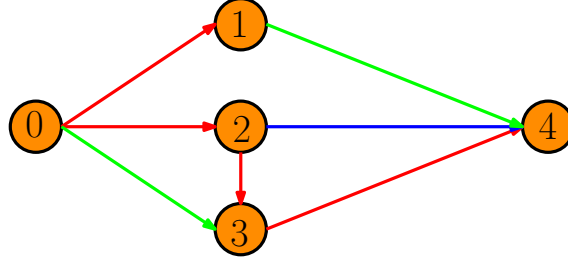


Figure 4.3: A simple graph, where edge color indicate the label. Let red be $\{a\}$, blue be $\{b\}$ and green be $\{c\}$. Vertex 0 can reach vertex 4 over several paths, but only has to store $\{a\}$, because all other paths are supersets of that path.

The maintenance algorithms below were built to work on an index with LI.

4.4.1 Adding an edge

Suppose we add an edge of the form (v, w, l) to our graph. Then any entry of the form (v', w', L) that visited v during construction might be altered. Only for vertex w we can be sure that no entry is modified, as any entry in $\text{Ind}(w)$ that could use (v, w, l) is already covered by an entry in $\text{Ind}(w)$.

Algorithm 12 shows how to run an update. Vertex v is pushed onto a queue q . While the queue is not empty, we pop a vertex x (initially v). We run an adapted version of Algorithm 3 called `LabelledBFSPerNode`, which will not attempt `TRYINSERT` if it can make a call to `FORWARDPROP`. This will not do a full reconstruction. Rather it will try to insert the entry (w, l) into $\text{Ind}(v)$ or take the entries from a landmark $v' \in V'$ that has been fully indexed.

4.4.2 Removing an edge

Removing an edge (v, w, l) consists of two steps. First we need to remove any entry $(w', L) \in \text{Ind}(v')$ for all $v' \in V$ which would not be in the index if (v, w, l) would not be an edge in the graph. Let A be the set of entries $(v', w', L) \in \text{Ind}(v')$ for all $v' \in V$ for which $v' \xrightarrow{L} w'$ before removing edge (v, w, l) and for which $v' \not\xrightarrow{L} w'$ after removing edge (v, w, l) . Next, we might need to rebuild parts of the index using `LabelledBFSPerNode` (Algorithm 3) for every $v \in V$, because entries that were not part of the index previously could be part of the index after removal of the edge.

The first difficulty lies in a good approximation of A . So far it appears we can only find some superset A' of A . The second difficulty has to do with the fact that we end up with an index that has some but not all of the entries. A call to `LabelledBFSPerNode` originated at a node w assumes that if an entry $(v, L) \in \text{Ind}(w)$ that we need not to look at any descendants of v . However, this might be wrong in case $(v, L) \in \text{Ind}(w)$, $(v, v', l) \in E$ and $(v, L \cup \{l\}) \notin \text{Ind}(w)$. The entry $(v, L \cup \{l\})$ might be removed by the first step.

Approximating A can be done by the following algorithm. After an edge (v, w, l) has been removed, we can (temporarily) promote v and w to a landmark and fully index these. Next, we iterate over all $w' \in V$. If there

Algorithm 12 $\text{insert}(v, w, l)$

```

1: add  $(v, w, l)$  to  $G$ 
2: reset hasBeenIndexed
3: hasBeenIndexed $(w) \leftarrow \text{True}$ 
4: let  $q$  be an empty queue
5: push  $v$  to  $q$ 
6: while  $q$  is not empty do
7:   if hasBeenIndexed $(w) = \text{True}$  then
8:     continue
9:   end if
10:   $x \leftarrow q.\text{pop}()$ 
11:  if  $x$  is a landmark then
12:     $\text{LabelledBFSPerNode}'(x)$ 
13:  end if
14:  for  $(v, x) \in E$  do
15:     $q.\text{push}(v)$ 
16:  end for
17:  hasBeenIndexed $(x) \leftarrow \text{True}$ 
18: end while

```

exists an entries $(v', L_1) \in \text{Ind}(w)$ and $(v', L_2) \in \text{Ind}(w')$ for some $v' \in V$. If $L_1 \cup \{l\} \subseteq L_2$ then we know that w' can reach v' with a L_2 -path and $l \in L_2$. Hence we can safely remove entry (v', L_2) from $\text{Ind}(w')$. Let A' be the resulting approximation and let A'_V be $\{v \mid (v, L) \in A'\}$.

We did not manage to find a 'good' solution for the second part. Either the solution we found for this part was too slow, i.e. not within 90% of the time of a full rebuild, or the solution would contain errors.

4.4.3 Changing edge label

As this is quite similar to first removing and then adding an edge we chose not to implement this case.

4.4.4 Adding a node

Adding a node without any edges is trivial and can be done in constant time.

4.4.5 Removing a node

As adding an edge or removing an edge is already quite expensive we do not expect there to be an efficient solution for removing a node. Hence we think a full rebuild is a better solution here.

4.5 Extensions

In this section we discuss several extensions that could be possible based on the indices we have discussed so far.

4.5.1 Query for all nodes

One of the many possible extensions could be to add support for the following type of query: "find all v' s.t. given a label set L and $v \in V$ we have that query (v, v', L) is true". The result of the query could be a mapping $m : V \rightarrow \{0, 1\}$. m has the same role as the variable **marked** in the description of BFS (Section 4.1.1) and we have that $m(v') = 1 \Leftrightarrow v \xrightarrow{L} v'$.

BFS could answer this query by simply continuing searching the graph rather than stopping at the target w . For each node $v \in V$ that is traversed, we set $m(v) \leftarrow 1$. This has a nice $O(|V| + |E|)$ running time. However this could be too much already.

LI could do the same as BFS, but make use of landmarked vertices. When a landmark v' is hit, we look at all (w, L') in $\text{Ind}(v')$ and test whether $L' \subseteq L$. If this is, we set $m(v) \leftarrow 1$.

4.5.2 Distance queries

Another extension could be to find the distance of a query as well, i.e. given a query $q = (v, w, L)$ find the shortest path P (w.r.t. the number of edges) that connects v and w such that $\mathbf{Labels}(P) \subseteq L$. This kind of query is answered by Bonchi et al. [1].

There is an important difference with “LCR” here though. We wish to find a distance d s.t. $d = \min_X(\mathbf{Len}(P))$ where $X = P \in \mathbf{Paths}(v, w) \wedge \mathbf{Labels}(P) \subseteq L$. This implies we need to store supersets for a pair v and w as well, given a superset has a shorter distance from v to w . Given two entries (v, w, L_1) and (v, w, L_2) connecting v and w with $L_1 \subset L_2$, we need to store both entries if the distance associated to the second entry is less than the distance of the first entry. This will increase the index size.

There are two things we need to change to make e.g. LI capable of answering these kinds of queries.

TRYINSERT needs to be modified s.t. it only adds an entry (v, w, L_1, d_1) if for all other entries (v, w, L_2, d_2) we either have that $L_1 \not\subseteq L_2 \wedge L_2 \subseteq L_1$ or that $d_1 < d_2 \wedge L_2 \subset L_1$. Moreover, it should remove all entries (v, w, L_2, d_2) when $L_1 \subseteq L_2$ and $d_2 \geq d_1$.

The entries for the heap in Algorithm 3 need to include a distance metric w.r.t. the number of edges as well. Preferably the heap should order its entries by this distance.

Chapter 5

Experimental design

In this chapter we describe the datasets and hardware that have been used in our experiments, as well as the way in which we generated queries for the experiments.

5.1 Datasets

For our experiments we used both synthetic and real data. Table 5.1 shows a full summary of all datasets including a number of statistics. In total 68 datasets were used.

5.1.1 Synthetic datasets

We generated the synthetic graphs using SNAP [7, 8] using either the ‘Preferential Attachment’ model (pa or PA) the ‘Erdos-Renyi’ model (er or ER), the ‘Forest Fire’ model (ff or FF) or the ‘Power Law’ model (pl or PL). In case a model (pa, er) would yield an undirected graph the direction would be randomly chosen. The edge labels can either have a uniform, normal ($\mu = |\mathcal{L}|/2, \sigma = |\mathcal{L}|/4$) or exponential ($\lambda = \frac{|\mathcal{L}|}{1.7}$) distribution. In the last case and in case $|\mathcal{L}| = 8$, we found that roughly 60% of the labels have the same value. A dataset with a uniform, normal or exponential label set distribution is also referred to as a uni-, norm- or exp-dataset.

The synthetic datasets are distinguishable according to five measures: the number of vertices (1k, 5k, 125k or 625k), the average degree per node (2 or 5), the number of labels (8, 10, 12, 14 or 16), the model (pa, ff, er, pl) and the label set distribution (uniform, normal or exponential).

5.1.2 Real datasets

The real datasets are taken from SNAP datasets [7] and KONECT¹ [6]. There are two types of real datasets: those which already had labels (r2) and those to which labels were added synthetically in the same manner as the synthetic datasets (r1). In the latter case, this is always an exponential distribution with 8 labels. In Figure 5.1, we can see the label distribution of some r2-datasets. We removed edges of the form $(v, v) \in E$ for any $v \in V$.

- **Robots** has been taken from Trustlet² and was published on June 2014. It is a trust-network where user A can give a certain level of trust to user B by using a certain type of edge.
- **Advogato** has been taken from KONECT and was published on April 2016. Advogato is an online community platform for developers of free software launched in 1999. Nodes are users of Advogato and the directed edges represent trust relationships. A trust link is called a “certification” on Advogato, and three different levels of certifications are possible on Advogato, corresponding to three different edge labels.
- **yagoFacts-small** is a sample taken from Max-Planck Institut³. YAGO is a huge semantic knowledge base, derived from Wikipedia WordNet and GeoNames. Currently, YAGO has knowledge of more than 10

¹<http://konect.uni-koblenz.de/networks/>

²<http://tinyurl.com/gnexfoy>

³<http://www.mpi-inf.mpg.de/departments/databases-and-information-systems/research/yago-naga/yago/>

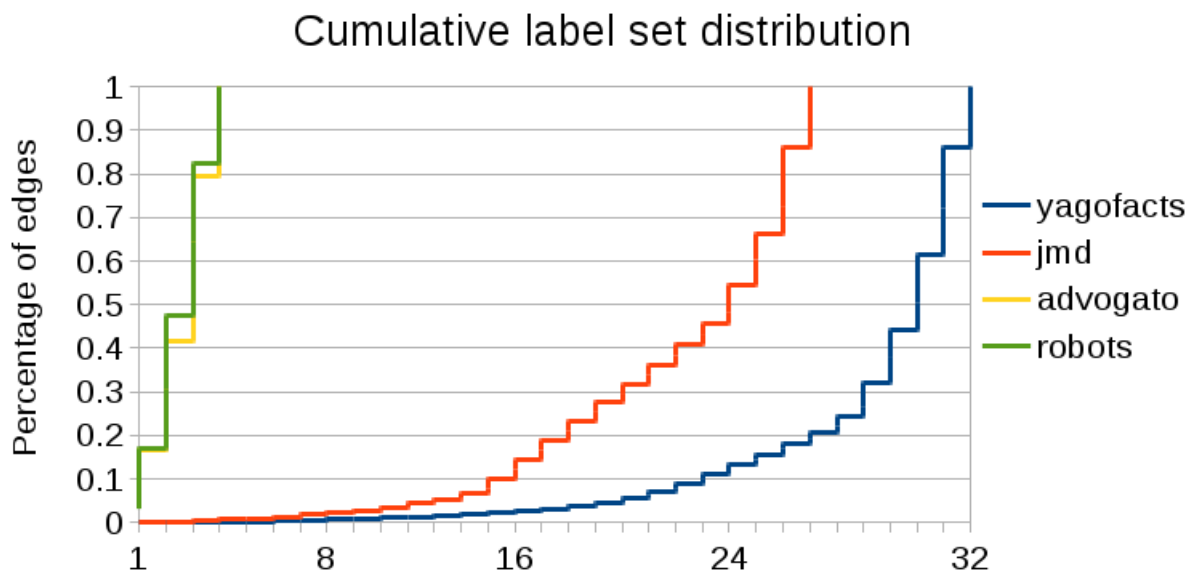


Figure 5.1: A cumulative distribution of the label set for the real datasets with an already existing label set distribution, i.e. the label set distribution was not synthetically added. **Robots** and **Advogato** have $|\mathcal{L}| = 4$, **jmd** has $|\mathcal{L}| = 25$ and **yagofacts** has $|\mathcal{L}| = 32$.

million entities (like persons, organizations, cities, etc.) and contains more than 120 million facts about these entities.

- **jmd** is a sample taken from the RDF-dump of Jamendo. Jamendo is a large repository of Creative Commons licensed music, based in France.
- **subeljCoraL8** has been taken from KONECT and was put online on April 2016. Nodes represent scientific papers. An edge between two nodes indicates that the left node cites the right node.
- **arXivheppL8exp** has been taken from KONECT and was put online on April 2016. It is the network of publications in the arXiv's High Energy Physics Phenomenology (hep-ph) section.
- **p2p-GnutellaL8exp** has been taken from KONECT. It's a network of Gnutella hosts from 2002. The nodes represent Gnutella hosts, and the directed edges represent connections between them. The dataset is from August 31, 2002.
- **socSlashdot0902L8exp** has been taken from SNAP. The website features user-submitted and editor-evaluated current primarily technology oriented news. The network contains friend/foe links between the users of Slashdot. The network was obtained in February 2009.
- **soc-sign-epinionsL8** has been taken from SNAP. It is a who-trust-whom online social network of a general consumer review site Epinions.com. Members of the site can decide whether to trust each other. The date it was obtained, is not given. We downloaded it on the 14th of April, 2016.
- **NotreDameL8exp** has been taken from KONECT. It is the directed network of hyperlinks between the web pages from the website of the University of Notre Dame. The date it was obtained, is not given. We downloaded it on the 14th of April, 2016.
- **webGoogle** has been taken from KONECT. This is a network of web pages connected by hyperlinks. The data was released in 2002 by Google as a part of the Google Programming Contest.
- **webBerkStan** has been taken from KONECT. This is the hyperlink network of the websites of the Universities in Berkley and Stanford. Nodes represent web pages, and directed edges represent hyperlinks.

- **webStanford** has been taken from KONECT. This is the directed network of hyperlinks between the web pages from the website of the Stanford University.
- **socPokecRelationShipsL8exp** has been taken from KONECT. This is the friendship network from the Slovak social network Pokec. Nodes are users of Pokec and directed edges represent friendships.
- **zhishihudongL8exp** has been taken from KONECT. These are “related to” links between articles of the Chinese online encyclopedia Hudong (<http://www.hudong.com/>).
- **usPatentsL8exp** has been taken from KONECT. This is the citation network of patents registered with the United States Patent and Trademark Office. Each node is a patent, and a directed edge represents a patent and an edge represents a citation.

5.1.3 Summary of datasets

Table 5.1 shows a full overview of all the datasets that were used in at least one experiment.

Table 5.1: A listing of all datasets with statistics. The model is one of the following: er (Erdos-Renyi), ff (Forest-Fire), pl (PowerLaw), pa (Preferential Attachment), r1 (real data, but synthetically added labels) and r2 (real data, real labels). max SCC % indicates the percentage of nodes that are in the largest strongly connected component (SCC). $\#\delta$ is the number of triangles in the graph, i.e. $u, v, w \in V$ s.t. $(u, v), (v, w), (w, u) \in E$. In the model column, pa stands for ‘preferential attachment’, pl for ‘powerlaw’ and ff for ‘forest fire’, whereas none indicates the dataset is a real one.

name	$ V $	$ E $	$ L $	model	max SCC %	$\#\delta$	dia
ERV1kD2L8uni	1,000	2,000	8	er	0.63	0	18
ERV1kD5L8uni	1,000	5,000	8	er	0.98	35	9
ff1k-0.2-0.4	1,000	2,552	8	ff	0.48	405	15
pl-1kL8a2.0exp	1,000	2,306	8	pl	0.42	200	10
plV1kL8a2.0exp	1,000	2,306	8	pl	0.41	195	8
V1kD2L12exp	1,000	1,997	12	pa	0.56	12	13
V1kD2L8exp	1,000	1,997	8	pa	0.61	15	16
V1kD2L8norm	1,000	1,997	8	pa	0.54	13	13
V1kD2L8uni	1,000	1,997	8	pa	0.55	16	13
V1kD5L8exp	1,000	4,985	8	pa	0.97	184	8
V1kD5L8norm	1,000	4,985	8	pa	0.97	212	8
V1kD5L8uni	1,000	4,985	8	pa	0.97	198	7
robots	1,484	2,960	4	r1	0.14	37	10
ERV5kD2L8uni	5,000	10,000	8	er	0.63	1	27
ff5k-0.2-0.4	5,000	12,718	8	ff	0.43	2,004	17
pl-5kL8a2.0exp	5,000	15,634	8	pl	0.44	2,246	10
plV5kL8a2.0exp	5,000	15,634	8	pl	0.44	2,237	10
V5kD10L8exp	5,000	49,945	8	pa	0.99	2,161	7
V5kD2L12exp	5,000	9,997	12	pa	0.55	31	16
V5kD2L8exp	5,000	9,997	8	pa	0.56	30	19
V5kD2L8norm	5,000	9,997	8	pa	0.56	27	20
V5kD2L8uni	5,000	9,997	8	pa	0.55	28	18
V5kD5L8exp	5,000	24,985	8	pa	0.97	406	9
advogato	5,417	51,327	4	r1	0.59	4,175	10
yagoFacts-small	10,000	22,122	32	r1	0.00	34	13
subeljCoraL8exp	23,167	91,500	8	r2	0.17	793	41
ERV25kD2L8uni	25,000	50,000	8	er	0.63	3	31
plV25L8ka2.0exp	25,000	90,449	8	pl	0.45	14,086	11
V25kD2L12exp	25,000	49,997	12	pa	0.55	43	23

Table 5.1: A listing of all datasets with statistics. The model is one of the following: er (Erdos-Renyi), ff (Forest-Fire), pl (PowerLaw), pa (Preferential Attachment), r1 (real data, but synthetically added labels) and r2 (real data, real labels). max SCC % indicates the percentage of nodes that are in the largest strongly connected component (SCC). $\#\delta$ is the number of triangles in the graph, i.e. $u, v, w \in V$ s.t. $(u, v), (v, w), (w, u) \in E$. In the model column, pa stands for 'preferential attachment', pl for 'powerlaw' and ff for 'forest fire', whereas none indicates the dataset is a real one.

name	$ V $	$ E $	$ L $	model	max SCC %	$\#\delta$	dia
V25kD2L8exp	25,000	49,997	8	pa	0.55	47	23
V25kD2L8norm	25,000	49,997	8	pa	0.55	52	24
V25kD2L8uni	25,000	49,997	8	pa	0.56	40	22
V25kD3L10exp	25,000	74,994	10	pa	0.84	154	15
V25kD3L12exp	25,000	74,994	12	pa	0.83	160	16
V25kD3L14exp	25,000	74,994	14	pa	0.83	155	16
V25kD3L16exp	25,000	74,994	16	pa	0.84	152	16
V25kD3L8exp	25,000	74,994	8	pa	0.83	159	17
V25kD4L10exp	25,000	99,990	10	pa	0.93	376	13
V25kD4L12exp	25,000	99,990	12	pa	0.93	421	14
V25kD4L14exp	25,000	99,990	14	pa	0.93	373	13
V25kD4L16exp	25,000	99,990	16	pa	0.93	388	13
V25kD4L8exp	25,000	99,990	8	pa	0.93	357	13
V25kD5L10exp	25,000	124,985	10	pa	0.97	666	11
V25kD5L12exp	25,000	124,985	12	pa	0.97	669	12
V25kD5L14exp	25,000	124,985	14	pa	0.97	653	12
V25kD5L16exp	25,000	124,985	16	pa	0.97	655	12
V25kD5L8exp	25,000	124,985	8	pa	0.97	645	11
arXivhepphL8exp	34,547	421,534	8	r2	0.36	248	45
p2p-GnutellaL8exp	62,587	147,892	8	r2	0.22	56	28
socSlashdotL8exp	82,168	870,161	8	r2	0.86	12,296	12
ERV125kD2L8uni	125,000	250,000	8	er	0.63	0	37
plV125L8ka2.0exp	125,000	518,207	8	pl	0.46	91,482	11
V125kD2L12exp	125,000	249,997	12	pa	0.56	71	24
V125kD2L8exp	125,000	249,997	8	pa	0.56	69	24
V125kD2L8norm	125,000	249,997	8	pa	0.56	81	27
V125kD2L8uni	125,000	249,997	8	pa	0.56	66	27
V125kD5L8exp	125,000	624,985	8	pa	0.97	994	14
soc-sign-epinionsL8exp	131,828	840,799	8	r2	0.31	73,928	17
webStanfordL8exp	281,904	2,312,497	8	r2	0.53	94,008	591
NotreDameL8exp	325,730	1,469,679	8	r2	0.16	52,100	72
citeseerL8exp	384,414	1,751,463	8	r2	0.04	1,566	71
twitterL8exp	465,018	834,797	8	r2	0.00	54	15
jmd	486,320	1,049,647	26	r1	0.00	0	5
V625kD5L8exp	625,000	3,124,985	8	pa	0.99	1894	16
plV625ka2.0L8exp	625,000	2,897,916	8	pa	0.97	350,125	14
webBerkstanL8exp	685,231	7,600,595	8	r2	0.48	318,370	670
webGoogleL8exp	875,713	5,105,039	8	r2	0.97	306,723	24
socPokecRelationshipsL8exp	1,632,803	30,622,534	8	r2	1	653,121	14
zhishihudongL8exp	2,452,715	18,854,882	8	r2	0.98	2,004,783	108
usPatentsL8exp	3,774,769	16,518,947	8	r2	0.00	0	23

5.2 Queries

We strived to generate a ‘realistic query workload’. By this we meant two things. Queries $q = (v, w, L)$ should not have a L -path P from v to w where $\text{Len}(P)$ is very low, because this favours BFS over other approaches, e.g. LI. We do not wish that a large part of our queries has the same starting point, as this might favour for instance LI.

For each dataset, both synthetic and real, we generated three query sets. The number of labels in a query varies between the three query sets and is: $\lfloor |\mathcal{L}|/4 \rfloor$, $\lfloor |\mathcal{L}|/2 \rfloor$ and $|\mathcal{L}| - 2$ for the datasets when $|\mathcal{L}| \geq 8$ and 1, 2 and 3 otherwise. Each query set consists of 200 (for graphs with $|E| < 5,000$) or 2,000 queries (otherwise). The first half consists of True-queries and the second half consists of False-queries.

For each query we took into account the difficulty. The difficulty of a query $q = (s, t, L)$ can be defined in the following way: the number of vertices visited during our implementation of BFS from s to t before reaching a conclusion, i.e. returning false or true. Our queue is a FIFO-queue and all edges for a given $v \in V$ of the form (v, w, l) are sorted in ascending order according to the vertex id of w . This difficulty can be influenced by five main factors: the distance between s and t in the graph, the number of labels in L , whether the result of the query is true or false, whether it is cyclic or whether it has a skewed out-degree distribution. False queries with a large label set may explore much larger sections of the graph than true queries with the same label set, because a True-query stops after hitting its target. Similarly a graph with a very skewed out-degree distribution has a low diameter and hence has more difficulty generating difficult queries. One can see that if we were to omit this requirement and the difficulty of all queries in a dataset would be for instance relatively low, it could benefit BFS’s performance compared to that of any index.

Queries for a particular query set are generated in the following way over a set of rounds. Let nq be the desired number of True- or False-queries. Each round starts by choosing a random vertex $v \in V$ and a random minimal difficulty d_{\min} . The minimal difficulty is between $\lceil \log_2(n) \rceil$ and $n/10$ (for $n \leq 500,000$) and between $\lceil \log_2(n) + 10 \rceil$ and $n/100$ otherwise. We generate up to $k' = nq/100$ random label sets L . The labels in the label set are chosen randomly according to a uniform distribution. Then, it is tested whether (v, w, L) is True or False using BFS and how many vertices were visited while doing this. Let d be this number of vertices. If $d \geq d_{\min}$, we add the query to the True- or False-set depending on its outcome until we have nq queries of both types. The parameters k' and d_{\min} can respectively increase and decrease by 1 for each 100 rounds passed.

The parameters mentioned are chosen to ensure that the queries either have no L -path or a L -path of sufficient length, that a small part of the queries have the same starting vertex $v \in V$ and that for any dataset we were able to generate queries in the end. For some datasets, particularly acyclic datasets like **jmd**, it is hard to generate difficult queries. This is why we allowed the parameters k' and d_{\min} to scale with the number of rounds. For a similar reason we chose to set the maximal difficulty to $n/100$ for $n \geq 500,000$, because the difficulty does not grow linearly with the number of nodes, at least for the True-queries.

5.3 Hardware

We used a server on the National Institute of Informatics (NII) in Tokyo, Japan. The server has 258GB of memory and a 2.9@Ghz 32-core processor. However, we did not let any of our experiments exceed an index size of 128GB and we set a 6 hour time limit. The experiments were all single-threaded. The OS is Linux Ubuntu.

5.4 Methods

Table 5.2 describes all methods used in at least one experiment with an abbreviation.

Table 5.2: Summary of all used methods with an abbreviation, description and main parameters.

method	abbreviation	description	parameters
BFS	BFS	Regular BFS, no indexing.	none
ZOU	Zou	Zou-method, exact index.	none
CLUSTEREDEXACT	Clus	Clustered-method, exact index.	K : number of clusters
NEIGHBOUREXCHANGE	Nei	Neighbour-exchange method, exact index.	none
PARTIALINDEX	Par	Partial-index, a partial index.	b : budget per node
LI	L	Landmarked, exact index, no extension.	k : number of landmarks (always required)
LI+EXTv1	L2	Landmarked second extension.	b : budget per non-landmark node
LI+OTH+EXTv1	L12	Landmarked first and second extension.	b : budget per non-landmark node
LI+OTH+EXTv2	L23	Landmarked second and third extension.	b : budget per non-landmark node.
DOUBLEBFS	DBFS	Landmarked where $k = N$.	None.

Chapter 6

Experiments

In this chapter we describe the experimental results of the methods described in Chapter 5, as summarized in Table 5.2.

6.1 Introduction

The experiments have been divided into three parts, based on the maximum number of edges that are in the datasets of each of these parts: $0 < |E| \leq 5,000$ (small), $5,000 < |E| \leq 500,000$ (medium) and $|E| \geq 500,000$ (large).

The first part includes all methods discussed in Chapter 5 (PARTIALINDEX, LI, ZOU, NEIGHBOUREXCHANGE, JOININDEX, CLUSTEREDEXACT and DOUBLEBFS) and is meant to select the more promising methods out of the less promising ones.

The second part shows the effect of the out-degree and label set distribution and the size of the label set $|\mathcal{L}|$ on the index construction time and index size.

The third part shows the limits of our major contribution, i.e. LI+OTH+EXTv2.

The goal of an index is to speed up the query answering process relative to BFS.

When we discuss the speed-ups over BFS in this section, we mean “total speed-ups”. A total speed-up over a query condition is equal to the sum of all query times of BFS over that query condition over the sum of all query times of a different method over that query condition. We do wish to note that the individual speed-ups can be very different from this. An “individual speed-up” is the time that a query q took using BFS over the time that it took using a different method. We may have that the “total speed-up” is about 2.0 but that about 30% of the queries has an “individual speed-up” of at least 100. An “average speed-up” is the average over all individual speed-ups for all queries within a certain query condition.

A query condition Q_c is a pair of the form $(\{True, False\}, |L|)$ where $L \subseteq \mathcal{L}$. Let n_q be the number of queries in a certain query condition and let q_j for $1 \leq j \leq n_q$ be the j ’th query belonging to Q_c . Let $T_M(Q_c, j)$ be the time taken by method M to answer query q_j belonging to Q_c . Table 6.1 shows the calculations for the speed-ups

Table 6.1: Speed-ups calculations, showing the difference between “total”, “average” and “individual” speed-up.

Name	Formula
Total speed-up for Q_c and M'	$\frac{\sum_{i=1}^{n_q} (T_{BFS}(Q_c, i))}{\sum_{i=1}^{n_q} (T_{M'}(Q_c, i))}$
Average speed-up for Q_c and M'	$\sum_{i=1}^{n_q} \frac{(T(Q_c, i))}{(T_{M'}(Q_c, i))}$
Individual speed-up for Q_c , q_j and M'	$T_{BFS}(Q_c, j)/T_{M'}(Q_c, j)$

6.2 Part 1: small graphs ($0 < |E| \leq 5,000$)

6.2.1 Datasets and methods

We used the methods PARTIALINDEX, LI+OTH($k = n/10, b = k/10$), LI+OTH+EXTv1($k = n/10, b = k/10$) ZOU, NEIGHBOUREXCHANGE, CLUSTEREDEXACT (with 5 clusters) and DOUBLEBFS. NEIGHBOUREXCHANGE used blocked mode, which improved the index construction time at the cost of a larger index size. k is the number of landmarks for LI and b is the number of entries $(v_i, L) \in v_0$ for non-landmark vertices $v_0 \in V$.

6.2.2 Index construction time (s) and size (MB)

Table 6.2: Index construction time (s) for the 7 methods and all datasets.

dataset	BFS	Zou	Clus	Nei	Par	L2	L12	DBFS
ERV1kD2L8uni	0.01	1,099.13	0.03	13.56	0.13	0.86	0.86	3.18
ERV1kD5L8uni	0.01	1,219.76	0.89	51.94	1.18	3.91	3.92	17.61
ff1k-0.2-0.4	0.01	43.48	0.06	3.38	0.14	0.06	0.06	0.22
plV1kL8a2.0exp	0.01	39.60	0.06	2.19	0.15	0.13	0.13	0.39
robots	0.01	11.99	0.07	1.49	0.07	0.06	0.06	0.11
V1kD2L12exp	0.01	39.96	0.03	3.15	0.07	0.03	0.03	0.19
V1kD2L8exp	0.01	41.48	0.03	3.51	0.06	0.03	0.03	0.20
V1kD2L8norm	0.01	308.75	0.04	6.22	0.12	0.17	0.17	0.55
V1kD2L8uni	0.01	502.83	0.05	6.04	0.15	0.25	0.25	0.75
V1kD5L8exp	0.01	191.39	0.17	9.74	0.38	0.28	0.27	1.01
V1kD5L8norm	0.01	611.15	0.30	17.53	0.91	0.53	0.53	2.60
V1kD5L8uni	0.01	908.33	0.37	19.15	1.01	0.69	0.69	3.44
Average	0.01	418.15	0.18	11.49	0.36	0.58	0.58	2.52

Table 6.3: Index size (MB) for the 7 methods and all datasets.

dataset	BFS	Zou	Clus	Nei	Par	L2	L12	DBFS
ERV1kD2L8uni	0.03	26.64	0.04	36.49	1.53	3.87	3.88	26.76
ERV1kD5L8uni	0.08	88.54	9.50	85.58	8.00	9.47	9.49	88.68
ff1k-0.2-0.4	0.04	8.61	1.03	15.52	1.70	1.52	1.54	8.73
plV1kL8a2.0exp	0.03	7.02	0.65	15.30	1.67	1.23	1.24	7.15
robots	0.04	4.84	0.90	28.26	1.15	1.76	1.78	4.95
V1kD2L12exp	0.03	10.70	0.37	15.78	1.41	1.36	1.39	10.87
V1kD2L8exp	0.03	11.66	0.38	15.87	1.48	1.47	1.48	11.77
V1kD2L8norm	0.03	17.35	0.47	23.74	1.90	2.61	2.63	17.47
V1kD2L8uni	0.03	21.36	0.77	27.96	2.26	3.33	3.35	21.48
V1kD5L8exp	0.07	29.26	4.23	26.10	4.65	3.84	3.86	29.39
V1kD5L8norm	0.07	49.21	4.82	46.40	7.69	5.86	5.89	49.35
V1kD5L8uni	0.07	57.28	5.75	54.40	8.88	6.73	6.76	57.42
Average	0.05	27.71	2.41	32.62	3.53	3.59	3.61	27.84

Tables 6.2 and 6.3 show the results for all datasets and methods. The abbreviations for the methods are listed in Table 4.2.

We see that ZOU needs far more time to build the index than the other methods. The index size of ZOU is almost identical to that of DOUBLEBFS. Especially in the cases of a uniform distribution we see that the construction time of ZOU explodes.

We took a look at the log files of the experiments. This shows us how much time each step of ZOU consumed. We see that ZOU needs a few seconds (e.g. 7 seconds for dataset **V1kD2L8uni**) for building the index of each SCC which is step 2 of the process. We know that checking the simpleness of a path P for each triple $(L(P), P, d)$ is computationally expensive (see Section 4.1.2).

The last steps, i.e. 7 to 9, took roughly 80% of the total time in most cases. During these steps we had to make some costly operations frequently, e.g. `Prune()` in combination with \odot . A lot of these operations are redundant, if we were to use some ordering of the vertices.

To give an example, suppose we have two inner vertices $v_{i,1}, v_{i,2} \in C_i$ where C_i is the i 'th SCC and $v_{i,1}$ has one out-edge $(v_{i,1}, v_{i,2}, l) \in E$. Moreover, suppose $v_{i,2}$ has already 'taken the entries' from every out-portal $p_i \in C_i$. For vertex $v_{i,1}$ it might be more efficient to just 'take the entries' of $v_{i,2}$ rather than doing this for all $p_i \in C_i$. However, this is not done.

We think that especially during these last steps a more efficient implementation should be possible. It should be possible to remove at least one step, as step 7 is not in the original paper. A more efficient implementation of the `Prune()`-operator thereby using a different data structure should be possible. However, as the first 6 already take considerably more time than `DOUBLEBFS` (which is on average 165 times faster), we do not think it can compete with the other methods. The gap between `DOUBLEBFS` on one side and ZOU on the other is simply too large. Hence we have decided to refrain from using ZOU in the next 2 parts of the experiment.

The index construction time of `NEIGHBOUREXCHANGE` is also high compared `DOUBLEBFS`. We were not able to come up with some form of pruning for `NEIGHBOUREXCHANGE`. Hence we have decided to refrain from using `NEIGHBOUREXCHANGE` in the next 2 parts of the experiment.

6.2.3 Total speed-ups

We chose to only analyse the total speed-ups of `CLUSTEREDEXACT`, `PARTIALINDEX` and `LI` (both approaches) and `DOUBLEBFS`, because `NEIGHBOUREXCHANGE` and ZOU have an index construction time that is not competitive with the other approaches to start with. Table 6.4 shows the total speed-up per query condition for `CLUSTEREDEXACT` and `PARTIALINDEX` and Table 6.5 does the same for `LI`.

The results for `CLUSTEREDEXACT` are bad. Any 0.01-value in the table indicates that the total speed-up of `CLUSTEREDEXACT` was at least 100 times slower in that case. The quality of the clustering seems to play a big role here. From the experiment logs we can see that a lot of direct and failed attempts were made to reach either an out-port or to reach the target from an in-port. This has to do with the relatively high number of edges between the clusters.

Looking at Tables 6.2 and 6.3 the index construction time and index size of `PARTIALINDEX` were good, but the total speed-ups achieved by this method are bad. Only for True-queries there is a significant speed-up. This is due to the fact that when we hit a landmark v' when answering query (v, w, L) we can either have a direct hit (i.e. there is a L -path from v' to w) or not. Unlike `LI` `PARTIALINDEX` does not have any advantage over `BFS` in the latter case.

When we look at Table 6.5 we see that the speed-ups for `LI+OTH+EXTv1` are slightly above those of `LI+OTH` in some cases and `LI+OTH` is clearly better in the remaining cases. There is not a clear winner here.

The speed-ups achieved by `DOUBLEBFS`, i.e. `LI` with $k = n$, are way higher than all the other speed-ups. However we must note here that the actual query answer times are very low for `LI+OTH+EXTv1` (many in the range of 10^{-7}). There might be some inaccuracy here. `LI+OTH+EXTv1` can answer any query immediately and does not have to initialize for instance a queue or bitset for any query. These might be explanations as to why the difference is so big.

The total speed-ups for ER-datasets and uniformly distributed datasets are much higher. This is due to the fact that ER-datasets have a more evenly distributed out-degree distribution and that uniformly distributed datasets allow for more minimal label sets between any pair of nodes.

We do wish to note that the speed-ups are total speed-ups. The average or individual speed-up can be very different. In the next section we look into this as well.

Table 6.4: Total speed-up per condition (query set and True/False-queries) for CLUSTERED_{EXACT} and PARTIAL_{INDEX}.

dataset	qs_1, true $ \mathcal{L} /4$	qs_1, false	qs_2, true $ \mathcal{L} /2$	qs_2, false	qs_3, true $ \mathcal{L} - 2$	qs_3, false
CLUSTERED _{EXACT}						
ff1k-0.2-0.4	0.11	0.01	0.08	0.01	0.33	0.01
ERV1kD2L8uni	0.01	0.01	0.01	0.01	0.01	0.01
plV1kL8a2.0exp	0.01	0.01	0.01	0.01	0.01	0.01
V1kD2L12exp	0.01	0.01	0.01	0.01	0.01	0.01
V1kD2L8exp	0.01	0.01	0.01	0.01	0.01	0.01
V1kD2L8norm	0.01	0.01	0.01	0.01	0.01	0.01
V1kD2L8uni	0.01	0.01	0.01	0.01	0.01	0.01
V1kD5L8exp	0.03	0.01	0.01	0.01	0.02	0.01
V1kD5L8norm	0.01	0.01	0.02	0.01	0.04	0.01
V1kD5L8uni	0.01	0.01	0.01	0.01	0.02	0.01
robots	0.01	0.01	0.01	0.01	0.01	0.01
	.02	.01	.01	.01	.04	.01
PARTIAL _{INDEX}						
ff1k-0.2-0.4	16.45	0.89	25.38	0.92	20.04	0.93
ERV1kD2L8uni	0.81	0.40	1.04	0.83	2.33	0.92
plV1kL8a2.0exp	5.90	0.93	9.43	0.95	16.61	0.95
V1kD2L12exp	2.85	0.92	7.31	0.94	32.55	0.96
V1kD2L8exp	5.44	0.93	15.97	0.98	24.93	1.08
V1kD2L8norm	2.05	0.82	3.10	0.89	3.56	0.92
V1kD2L8uni	1.40	0.77	3.87	0.85	11.65	0.91
V1kD5L8exp	18.97	0.89	16.72	0.80	40.05	0.85
V1kD5L8norm	10.54	0.73	18.90	0.82	17.79	0.85
V1kD5L8uni	5.83	0.75	19.08	0.85	23.21	0.87
robots	2.30	0.95	4.04	0.93	7.49	0.94
	6.59	.81	11.34	.88	18.20	.92

Table 6.5: Total speed-up per condition (query set and True/False-queries) for LI+OTH ($k = n/10$) and ($k = n$) and DOUBLEBFS.

method name	qs_1, true $ \mathcal{L} /4$	qs_1, false	qs_2, true $ \mathcal{L} /2$	qs_2, false	qs_3, true $ \mathcal{L} - 2$	qs_3, false
LI+OTH	(k=n/10,b=k/10)					
ff1k-0.2-0.4	50.35	67.78	82.62	63.69	78.95	72.48
ERV1kD2L8uni	1.49	1.45	9.39	5.83	23.70	10.23
plV1kL8a2.0exp	134.04	158.41	143.64	134.72	170.08	162.39
V1kD2L12exp	51.34	97.17	42.83	102.12	88.46	101.84
V1kD2L8exp	66.42	122.55	66.05	138.35	73.42	172.11
V1kD2L8norm	15.59	10.48	30.23	33.30	50.58	85.88
V1kD2L8uni	12.34	6.43	24.31	34.37	43.75	91.44
V1kD5L8exp	31.20	2.91	37.13	28.90	75.18	16.54
V1kD5L8norm	22.02	31.86	28.73	11.47	30.88	4.23
V1kD5L8uni	23.01	34.33	31.13	24.05	38.09	3.57
robots	63.67	35.84	126.30	112.40	137.61	244.54
	42.86	51.74	56.57	62.65	73.70	87.75
LI+OTH+EXTv1	(k=n/10,b=k/10)					
ff1k-0.2-0.4	54.69	51.79	85.21	51.80	73.40	51.90
ERV1kD2L8uni	1.51	1.47	9.37	5.94	24.00	10.96
plV1kL8a2.0exp	139.65	134.77	144.60	111.14	168.89	118.95
V1kD2L12exp	63.32	85.99	76.97	90.78	92.60	86.21
V1kD2L8exp	67.47	85.15	65.09	130.52	71.35	82.27
V1kD2L8norm	14.82	9.65	30.20	30.23	47.05	66.67
V1kD2L8uni	11.45	6.33	23.59	33.85	48.85	79.25
V1kD5L8exp	31.14	3.68	36.84	28.86	73.24	20.78
V1kD5L8norm	23.50	30.62	32.05	18.38	32.72	6.94
V1kD5L8uni	23.25	33.88	31.18	28.62	39.71	4.91
robots	61.14	27.07	117.69	85.42	133.44	250.62
	44.72	42.76	59.34	55.95	73.20	70.86
DOUBLEBFS						
ff1k-0.2-0.4	155.91	2,844.84	178.56	3,220.82	199.93	3,537.62
ERV1kD2L8uni	7.17	64.35	41.55	416.49	211.83	1,953.72
plV1kL8a2.0exp	307.42	2,745.11	334.78	2,577.43	359.59	3,153.61
V1kD2L12exp	163.63	2,538.65	205.62	3,140.89	239.10	3,315.53
V1kD2L8exp	142.05	2,536.86	136.78	3,059.71	192.43	3,907.76
V1kD2L8norm	45.45	413.86	106.00	1,209.13	159.92	2,418.49
V1kD2L8uni	39.08	328.03	101.76	1,243.88	128.82	2,453.83
V1kD5L8exp	193.94	5,292.36	180.60	1,628.22	250.92	3,836.84
V1kD5L8norm	143.23	2,081.43	181.40	5,600.57	175.60	7,547.58
V1kD5L8uni	114.19	2,311.32	180.97	6,473.33	256.51	8,372.60
robots	97.79	695.46	260.54	1,889.16	304.76	2,379.20
	128.16	1,986.57	173.50	2,769.05	225.40	3,897.88

6.3 Part 2: medium graphs ($5,000 < |E| \leq 500,000$)

6.3.1 Datasets and methods

In this experiment we included all datasets that have more than 5,000 edges and at most 500,000 edges. The methods we used include LI+OTH ($k = n/20, b = 20$), LI+OTH ($k = n/20, b = 0$), LI+OTH+EXTv1 ($k = n/20, b = 20$) and LI+OTH+EXTv1 ($k = n/10, b = 20$). We set the budget to either 20 or 0. In case $b = 0$ one could argue that LI+OTH is the same as LI. k is the number of landmarks and b the budget per non-landmark.

For the datasets of the type Preferential-Attachment (or PA-datasets) that have 25,000 vertices and with a degree ≥ 3 and $L > 8$ we only ran LI+OTH+EXTv1 ($k = n/10, b = 20$)

6.3.2 Index construction time (s)

Table 6.6: Index construction time (s) for the 4 methods and all datasets. For the dataset **ERV125kD2L8uni** no index was produced within the 6 hours time limit. The table has been divided into a few sections: the top one contains the real datasets, the second one the ER-datasets, the third one the Forest-Fire datasets, the fourth one the Power-Law datasets and the last one the Preferential-Attachment datasets.

dataset	BFS	L12($\frac{n}{20}, 20$)	L1($\frac{n}{20}, 20$)	L2($\frac{n}{20}, 0$)	L12($\frac{n}{10}, 20$)
advogato	0.01	0.84	0.82	0.80	1.48
yagoFacts-small	27.87	28.04	28.04	28.03	28.04
subeljCoraL8exp	0.01	5.27	5.16	5.09	8.94
arXivhepphL8exp	0.06	291.75	292.83	293.22	417.95
p2p-GnutellaL8exp	0.02	103.05	102.19	101.50	173.48
ERV5kD2L8uni	0.01	26.98	27.03	26.98	41.33
ERV25kD2L8uni	0.02	978.67	977.01	1,069.33	1,448.87
ERV125kD2L8uni	-	-	-	-	-
ff5k-0.2-0.4	0.01	1.05	1.03	1.02	1.58
plV5kL8a2.0exp	0.01	1.48	1.46	1.45	2.40
plV25L8ka2.0exp	0.03	30.44	30.37	31.08	50.17
V5kD2L8exp	0.01	0.36	0.35	0.34	0.61
V5kD2L8norm	0.01	3.49	3.46	3.45	5.09
V5kD2L8uni	0.01	5.51	5.55	5.48	7.78
V25kD2L8exp	0.01	9.81	9.65	9.48	16.54
V25kD2L8norm	0.02	110.03	110.09	109.99	156.99
V25kD2L8uni	0.01	157.78	157.49	157.38	216.81
V125kD2L8exp	0.03	262.30	267.49	254.91	445.85
V125kD2L8norm	0.03	2,944.69	2,956.99	2,953.46	4,182.90
V125kD2L8uni	0.03	4,341.67	4,380.90	4,372.79	6,080.10
Total average	1.48	489.64	492.52	496.09	699.31

The top section of Table 6.6 contains the real datasets. There is quite some difference when we look at the index construction times of these datasets, e.g. **subeljCoraL8exp** (23k edges) took roughly 5 seconds whereas **arXivhepphL8exp** (34k edges) took way longer. A clear difference between these two is the out-degree distribution. **subeljCoraL8exp** has a more skewed out-degree distribution, i.e. a few nodes have a very high degree, whereas **arXivhepphL8exp** has a more balanced out-degree distribution. We have seen this effect multiple times, e.g. the ER-datasets have an almost uniform out-degree distribution.

The graph construction time of **yagoFacts-small** is quite high: 27.87 seconds. The only thing that makes this dataset very different from other datasets is the fact that it has $|\mathcal{L}| = 32$. However this does not explain why the graph construction time is this high.

The ER-datasets have a uniform label set distribution and a close to uniform out-degree distribution. Hence these datasets are difficult to process. This explains why we were unable to build an index within the time limit of 6 hours (21,600 s) for the dataset **ERV125kD2L8uni**. From **ERV125kD2L8uni**'s experiment logs we discovered that the index construction process had built around 50% of the landmarks, when the time limit was reached. By that time the index had reached a size of 20GB.

Looking at the PA-datasets and the label set distribution in Table 6.6 we see a large difference between the datasets with a uniform and normal label set distribution on the one hand and the exponential distribution on the other hand. For $|V| = 5k$, $|V| = 25k$ and $|V| = 125k$ we can see that the index construction takes about 13 times longer for a dataset with a uniform label set distribution than for that same dataset with an exponential distribution.

6.3.3 Index size (MB)

Table 6.7: Index size (MB) for the 4 methods and all datasets. The sections in the table are the in Table 6.6.

dataset	BFS	L12($\frac{n}{20}, 20$)	L1($\frac{n}{20}, 20$)	L2($\frac{n}{20}, 0$)	L12($\frac{n}{10}, 20$)
advogato	0.82	27.57	27.46	27.44	54.60
yagoFacts-small	0.28	0.51	0.30	0.30	0.74
subeljCoraL8exp	1.46	54.94	54.67	54.67	109.86
arXivhepphL8exp	6.74	613.21	611.96	611.92	1,281.71
p2p-GnutellaL8exp	2.36	3,668.12	3,667.25	3,667.24	7,298.90
ERV5kD2L8uni	0.16	62.00	61.97	61.97	119.71
ERV25kD2L8uni	0.80	1,645.72	1,645.56	1,645.56	3,151.02
ERV125kD2L8uni	-	-	-	-	-
ff5k-0.2-0.4	0.20	17.08	17.03	16.93	32.22
plV5kL8a2.0exp	0.25	17.38	17.29	17.29	32.71
plV25L8ka2.0exp	1.44	421.25	420.68	420.68	811.87
V5kD2L8exp	0.15	18.12	18.08	17.99	35.22
V5kD2L8norm	0.15	39.09	39.05	38.94	71.66
V5kD2L8uni	0.15	46.39	46.35	46.25	84.70
V25kD2L8exp	0.79	437.62	437.38	437.16	857.19
V25kD2L8norm	0.79	992.47	992.25	992.00	1,830.24
V25kD2L8uni	0.79	1,169.63	1,169.43	1,169.17	2,121.92
V125kD2L8exp	3.99	10,787.30	10,786.07	10,785.39	21,291.71
V125kD2L8norm	3.99	24,381.91	24,380.83	24,380.22	45,072.80
V125kD2L8uni	3.99	29,044.04	29,043.01	29,042.46	53,146.00
Total average	1.54	3,865.49	3,865.09	3,864.93	7,231.83

The top section of Table 6.7 shows the real datasets. An interesting observation here is that construction time is not directly related to index size. For example **p2p-GnutellaL8exp** has the largest index size in the last column 7.3GB roughly, which took less than 2 minutes to build. In contrast, **arXivhepphL8exp** took almost 5 minutes and ended up with a smaller index. An explanation for this could be in the fact that **p2p-GnutellaL8exp** has a more skewed degree distributions than **arXivhepphL8exp**. Landmarks are those vertices with a high total degree. When $j < k$ landmarks have already been constructed and a landmark $v_j \in V$ needs to be constructed, there is a higher chance in a graph like **p2p-GnutellaL8exp** (with a skewed out-degree distribution) that it will hit a landmark v_i with $0 < i \leq j$ than in a graph like **arXivhepphL8exp**.

There is a high rate of growth of the index size, looking at the second section of Table 6.7 which contains the ER-datasets. In the last column of the table, we see that the index size grows from 119MB to 3,151MB (26 increase). Dataset **ERV25kD2L8uni** has a larger index size (3.1GB) compared to its PA-counterpart (2.1GB)

looking at the last column.

There is a high rate of growth of the index size, looking at the fifth section of Table 6.7 which contains the PA-datasets. The index size of a PA-exp-dataset (i.e. with an exponential label set distribution) grows in the following way: $11.77 \rightarrow 35.22 \rightarrow 857.19 \rightarrow 21,291.92$. This makes for a rate of growth of roughly 24, except for the first step. In case of a PA-norm-dataset the numbers are: $17.47 \rightarrow 71.66 \rightarrow 1,830.24 \rightarrow 45,072.80$. This makes for a rate of growth of roughly 25, except for the first step. In case of a PA-uni-dataset we get: $21.48 \rightarrow 84.7 \rightarrow 1,830.24 \rightarrow 53,146$. This makes for an average rate of growth of 27 for the last two steps. This demonstrates that each time the number of vertices is multiplied by a factor 5, the index size of a PA-dataset indiscrimination of the label set distribution is multiplied by at least a factor 24.

However, there is a clear difference between the datasets with a uniform and normal label set distribution on the one hand and the exponential distribution on the other hand. If we look at the last column again of the fifth section of Table 6.7, a PA-uni-dataset has an index size roughly 2.4 times larger than its corresponding PA-exp-dataset for $|V| = 5k$, $|V| = 25k$ and $|V| = 125k$. For example, **V25kD2L8uni** has an index size of 2,122MB and **V25kD2L8exp** has one of 857MB.

6.3.4 Index construction time (s) and index size (MB) when $|\mathcal{L}| \geq 8$

Table 6.8: Index construction time (s) for the datasets of the type 'Preferential-Attachment' having 25,000 edges and a degree of either 3, 4 or 5 and at least 8 labels using LI+OTH+EXTv1 ($k = N/20, b = 20$).

Degree	$ \mathcal{L} = 8$	$ \mathcal{L} = 10$	$ \mathcal{L} = 12$	$ \mathcal{L} = 14$	$ \mathcal{L} = 16$
3	33.75	70.92	223.31	1,058.74	4,219.61
4	72.44	215.71	721.78	3,554.97	-
5	113.46	353.74	1,414.50	6,527.27	-

Table 6.9: Index size (MB) for the datasets of the type 'Preferential-Attachment' having 25,000 edges and a degree of either 3, 4 or 5 and at least 8 labels using LI+OTH+EXTv1 ($k = N/20, b = 20$).

Degree	$ \mathcal{L} = 8$	$ \mathcal{L} = 10$	$ \mathcal{L} = 12$	$ \mathcal{L} = 14$	$ \mathcal{L} = 16$
3	736.01	953.55	1,496.42	2,783.56	4,219.61
4	1,022.86	1,566.40	2,677.61	5,404.49	-
5	1,250.53	2,004.95	3,720.64	7,238.79	-

We had a different set-up for the PA-exp-datasets with $n \geq 25,000$ and $D \geq 3$ and $|\mathcal{L}| \geq 8$. Table 6.8 and 6.9 show the index construction time (s) and index size (MB) for these datasets. Figures 6.1 and fig:papfig2 show the same data. When $D \geq 4$ and $|\mathcal{L}| = 16$, we were unable to build an index within the 6 hours time limit.

There is an interaction effect between the degree D and the size of the label set $|\mathcal{L}|$ looking at Figures 6.1 and fig:papfig2. The line with $D = 5$ has the strongest growth which can be seen in the figures. Increasing the label set size by 2 will make the index construction time and index size grow exponentially. This growth is much stronger for the index construction time. This can be explained by the fact that the index construction time complexity of LI+OTH+EXTv1 has a $2^{|\mathcal{L}|}$ -term and a $m = D \cdot n$ -term, whereas the time complexity the index size does not have a m -term.

6.3.5 Total speed-ups

Table 6.10 shows the mean query answering time (ms) of BFS. We can see that the mean query answering time is often a low value in an order of magnitude of 10^{-4} seconds.

In Tables 6.11, 6.10 and 6.12 we see the total speed-up for queries of the first, second and third query sets respectively. We can compare the different methods against each other by looking at the respective columns.

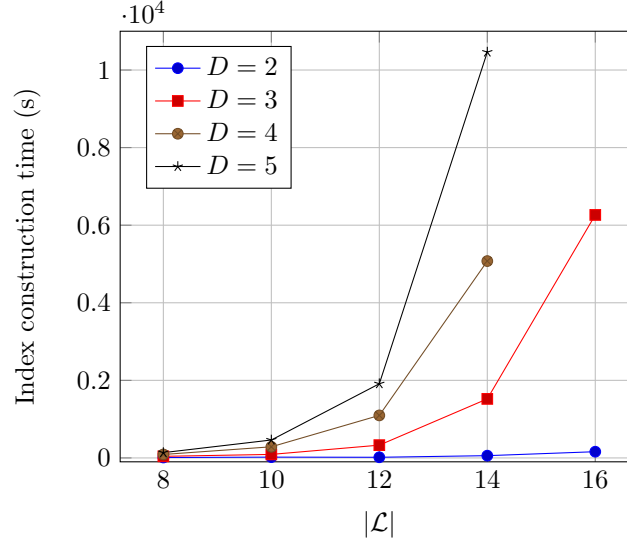


Figure 6.1: Index construction time (s) for PA-datasets with $n = 25,000$, as a function of the label set size $|\mathcal{L}|$. The different lines indicate the degree (either 2,3,4 or 5) of the datasets.

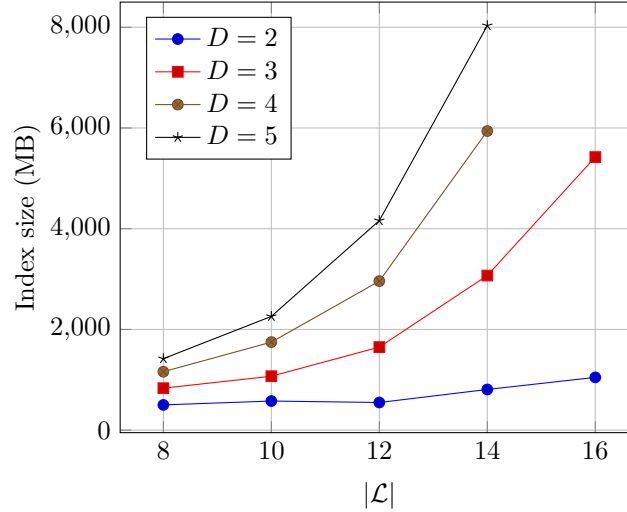


Figure 6.2: Index size (MB) for PA-datasets with $n = 25,000$, as a function of the label set size $|\mathcal{L}|$. The different lines indicate the degree (either 2,3,4 or 5) of the datasets.

When we look at the first two columns and compare these against the third and fourth column in the three tables, we do not observe any advantage of using the first extension. Only in a few cases there is a small improvement, e.g. **advogato** in Table 6.11. In the experiments leading up to this experiment we found some gain in using the first extension, but this gain was only exhibited for larger graphs. Therefore we will keep using the first extension.

There is a clear difference between LI+OTH with $b = 0$ and $b = 20$ looking at the three tables. The True-queries have a clear improvement in the $b = 20$ -case over the $b = 0$ -case. For the False-queries there is little improvement or even some disimprovement.

Finding the ideal budget size remains an open issue. We can imagine that a very large budget could increase the query answering time. For instance, given a query (v, w, L) we might have that v and w are in close proximity of each other, whereas any landmark $v' \in V'$ is not so close. In this case, a high b -value might first visit a large number landmarks before finding w through graph exploration.

Table 6.10: Mean query execution time (ms) of BFS for all datasets per query condition.

method name	qs_1, true $ \mathcal{L} /4$	qs_2, true $ \mathcal{L} /2$	qs_3, true $ \mathcal{L} - 2$	qs_1, false	qs_2, false	qs_3, false
advogato	0.203	0.243	0.246	0.297	0.383	0.370
yagoFacts-small	0.001	0.001	0.002	0.000	0.004	0.004
subeljCoraL8exp	0.225	0.240	0.261	0.337	0.378	0.420
arXivhepphL8exp	1.305	1.077	0.937	1.758	2.139	2.606
p2p-GnutellaL8exp	0.427	0.297	0.302	1.965	2.185	2.241
ERV5kD2L8uni	0.000	0.008	0.079	0.000	0.013	0.121
ERV25kD2L8uni	0.001	0.052	0.389	0.002	0.058	0.674
ff5k-0.2-0.4	0.067	0.043	0.044	0.183	0.125	0.169
plV5kL8a2.0exp	0.146	0.155	0.160	0.134	0.173	0.220
plV25L8ka2.0exp	0.841	0.862	0.917	0.823	0.929	1.341
V5kD2L8exp	0.059	0.048	0.049	0.154	0.177	0.191
V5kD2L8norm	0.018	0.042	0.037	0.025	0.076	0.139
V5kD2L8uni	0.007	0.029	0.033	0.007	0.063	0.135
V25kD2L8exp	0.207	0.170	0.147	0.807	0.893	0.877
V25kD2L8norm	0.303	0.143	0.139	0.195	0.385	0.738
V25kD2L8uni	0.041	0.141	0.115	0.049	0.367	0.749
V125kD2L8exp	0.304	0.260	0.279	4.243	3.901	5.066
V125kD2L8norm	0.235	0.267	0.253	0.402	1.856	3.786
V125kD2L8uni	0.101	0.264	0.163	0.251	1.662	4.008
Average	0.24	0.23	0.24	0.61	0.83	1.26

Increasing the number of landmarks, i.e. going from $\frac{N}{20}$ to $\frac{N}{10}$, improves the results to some extent. The total speed-up of the True-queries gets roughly doubled. The total speed-up of the False-queries does not always improve.

There is an effect of the label set distribution on the speed-ups. When we look at the total speed-ups for datasets that have an exponential label set distribution, e.g. **V5kD2L8exp**, and compare these against the speed-ups for their normal or uniform counterparts, e.g. **V5kD2L8uni** and **V5kD2L8norm**, we see that the speed-ups in the first case are higher. The datasets with an exponential label set distribution have one label l that is on the majority of the edges. The presence of l in a True-query (v, w, L) , i.e. $l \in L$, increases the probability that v can reach several landmarks $v' \in V$ and that one of these landmarks has a L -path to w .

There is an effect of the out-degree distribution on the speed-ups. The speed-ups for PL-datasets and **subeljCoraL8exp** and **p2p-GnutellaL8exp** are relatively high. All these datasets have a skewed out-degree distribution which favours the landmarked approach.

6.3.6 Total speed-ups when $|\mathcal{L}| \geq 8$

Table 6.14 shows the speed-ups for the PA-exp-datasets with $|V| \geq 25,000$, $D \geq 3$ and $|\mathcal{L}| \geq 8$.

The total speed-ups seem to increase with the degree for both True- and False-queries and seem to decrease with the label set size for the False-queries, but the numbers do not give a consistent picture.

6.3.7 Individual speed-ups

The speed-ups discussed in the previous section were “total speed-ups”. One should be careful here. The actual speed-up per query could be very different.

In Figure 6.3 we can see the speed-up per query for **p2p-GnutellaL8exp** and the third query set, i.e. query conditions $(\text{True}, |\mathcal{L}| - 2)$ and $(\text{False}, |\mathcal{L}| - 2)$ using method LI+OTH+EXTv1 ($k = \frac{n}{10}, b = 20$). The total

Table 6.11: Total speed-up for query conditions (True, $|\mathcal{L}|/4$) and (False, $|\mathcal{L}|/4$).

dataset	LI+OTH+EXTv1		LI+OTH		LI+OTH+EXTv1		LI+OTH+EXTv1	
	(20, $n/20$)		(20, $n/20$)		(0, $n/20$)		(20, $n/10$)	
ff5k-0.2-0.4	307.78	185.77	311.90	288.89	134.57	293.73	286.96	345.31
ERV5kD2L8uni	0.99	1.06	1.00	1.11	1.01	1.26	1.09	1.16
ERV25kD2L8uni	2.87	3.14	2.98	3.11	3.04	3.56	3.44	3.39
plV5kL8a2.0exp	651.29	319.53	661.20	375.38	346.68	423.30	701.33	662.25
plV25L8ka2.0exp	2,309.31	1,279.65	2,425.09	1,441.89	1,253.13	1,549.01	2,691.57	2,127.51
V5kD2L8exp	273.93	215.95	278.27	294.13	155.54	319.45	315.30	503.51
V5kD2L8norm	73.07	73.16	71.44	76.42	32.25	76.42	72.68	78.02
V5kD2L8uni	23.78	21.94	25.13	24.35	15.00	25.31	27.66	24.64
V5kD5L8exp	259.50	66.00	250.83	78.27	126.80	85.98	224.72	80.78
V25kD2L8exp	699.84	978.21	706.24	1,223.47	347.02	1,240.32	741.49	1,392.55
V25kD2L8norm	783.50	314.64	784.16	356.63	392.31	311.33	811.60	342.90
V25kD2L8uni	112.71	93.06	111.74	94.05	63.07	92.97	122.59	97.22
V125kD2L8exp	416.95	3,104.94	417.12	4,100.90	286.31	4,150.62	434.50	4,679.52
V125kD2L8norm	323.78	451.99	338.66	464.71	206.92	423.22	341.30	474.97
V125kD2L8uni	148.44	263.33	143.24	277.81	93.24	255.58	165.21	288.98
advogato	586.82	4.74	584.11	4.27	375.09	4.27	719.30	28.33
yagoFacts-small	4.55	1.83	4.52	2.08	3.32	2.13	5.27	1.95
subeljCoraL8exp	615.57	234.99	631.11	238.86	311.46	247.32	681.96	386.32
arXivhepphL8exp	687.32	4.97	706.53	4.84	627.12	4.83	1,183.89	20.56
p2p-GnutellaL8exp	1,027.91	5.25	1,044.45	5.49	576.22	5.47	1,346.45	2,020.55
Average	452.63	466.02	460.16	586.53	260.78	595.12	525.91	771.16

speed-ups are 771.22 and 646.22 (Table 6.13). About 55% of the True-queries and 93% of the False-queries have an individual speed-up that is above these values.

Table 6.12: Total speed-up for query conditions (True, $|\mathcal{L}|/2$) and (False, $|\mathcal{L}|/2$).

dataset	LI+OTH+EXTv1		LI+OTH		LI+OTH+EXTv1		LI+OTH+EXTv1	
	(20, $n/20$)		(20, $n/20$)		(0, $n/20$)		(20, $n/10$)	
ff5k-0.2-0.4	164.73	150.02	176.95	216.47	76.16	232.80	196.63	370.89
ERV5kD2L8uni	11.23	12.74	11.47	13.16	10.38	13.83	17.66	24.97
ERV25kD2L8uni	55.11	37.42	55.14	37.46	47.56	39.47	84.21	72.92
plV5kL8a2.0exp	757.26	349.79	756.66	558.84	395.94	591.28	818.06	516.40
plV25L8ka2.0exp	2,491.84	1,239.38	2,598.33	1,501.20	1,342.06	1,586.26	3,074.60	1,989.90
V5kD2L8exp	217.68	239.92	226.79	360.22	109.30	376.70	225.93	342.46
V5kD2L8norm	156.69	166.32	164.23	172.24	63.92	183.34	173.58	190.85
V5kD2L8uni	131.64	160.49	127.38	167.17	76.06	190.87	123.10	179.01
V5kD5L8exp	150.14	22.46	157.36	21.59	81.10	22.73	192.73	182.00
V25kD2L8exp	534.99	660.86	532.77	1,005.91	256.72	1,030.47	589.99	1,378.69
V25kD2L8norm	553.15	626.96	541.12	666.18	300.78	545.65	553.74	661.91
V25kD2L8uni	353.07	554.60	349.54	598.65	183.86	590.92	389.71	646.90
V125kD2L8exp	395.02	1,753.97	397.41	3,519.16	265.31	3,585.02	410.49	3,129.41
V125kD2L8norm	352.67	1,851.65	358.91	1,909.63	225.81	1,755.74	362.25	2,004.46
V125kD2L8uni	324.63	1,488.86	328.73	1,584.35	213.34	1,439.30	409.34	1,616.89
advogato	696.33	3.39	689.43	2.89	483.19	2.89	976.80	14.66
yagoFacts-small	3.86	9.23	4.08	10.84	2.52	10.65	4.73	9.59
subeljCoraL8exp	703.14	238.04	728.82	249.90	312.11	249.98	825.56	453.70
arXivhepphL8exp	959.57	5.45	981.89	5.01	610.50	5.00	1,662.81	21.73
p2p-GnutellaL8exp	668.38	2.73	672.63	2.94	401.28	2.95	771.22	646.25
Average	459.59	578.57	468.49	707.27	255.99	703.59	552.39	832.52

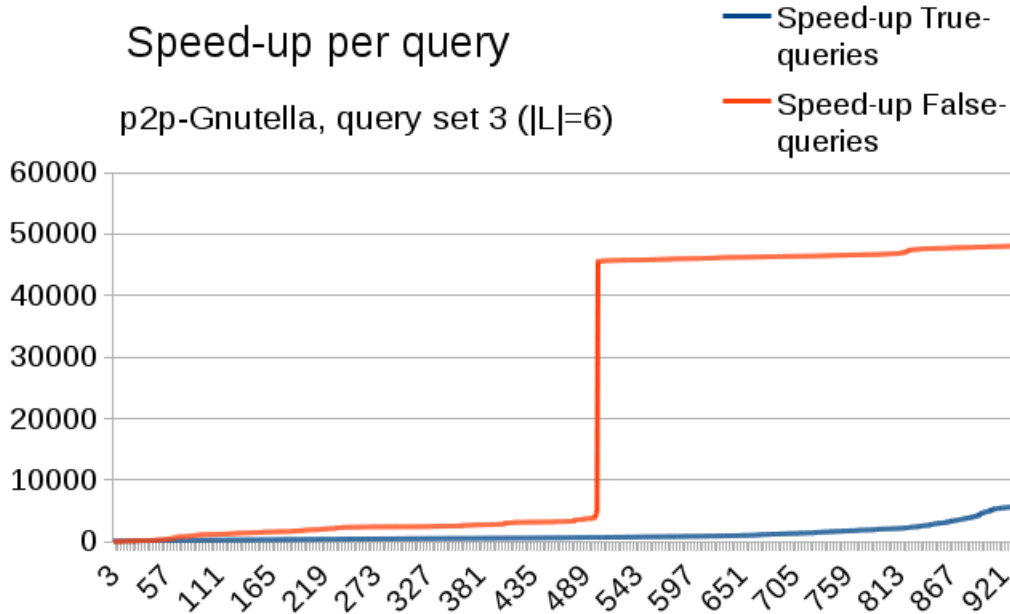
Figure 6.3: Individual speed-ups (sorted in ascending order) for **p2p-GnutellaL8exp** and query conditions (True, $|\mathcal{L}| - 2$) and (False, $|\mathcal{L}| - 2$) using LI+OTH+EXTv1 ($k = n/10, b = 20$).

Table 6.13: Total speed-up for query conditions (True, $|\mathcal{L}| - 2$) and (False, $|\mathcal{L}| - 2$).

dataset	LI+OTH+EXTv1		LI+OTH		LI+OTH+EXTv1		LI+OTH+EXTv1	
	(20, $n/20$)		(20, $n/20$)		(0, $n/20$)		(20, $n/10$)	
ff5k-0.2-0.4	217.50	158.84	202.46	277.44	104.16	292.53	194.98	390.77
ERV5kD2L8uni	163.49	5.46	166.75	5.29	106.43	5.34	225.74	76.79
ERV25kD2L8uni	610.40	4.72	604.85	4.72	393.05	4.73	932.80	389.21
plV5kL8a2.0exp	781.53	308.26	786.45	464.32	413.91	498.58	779.24	470.15
plV25L8ka2.0exp	2,694.00	1,137.77	2,832.51	1,649.38	1,361.76	1,673.05	2,846.10	2,020.66
V5kD2L8exp	232.45	225.10	246.24	234.68	113.71	247.78	252.81	414.74
V5kD2L8norm	164.18	192.42	169.53	212.75	62.33	210.08	167.39	302.03
V5kD2L8uni	139.82	326.93	147.65	324.28	67.04	370.80	142.46	361.71
V5kD5L8exp	180.94	278.20	169.91	273.02	74.39	314.76	200.15	399.34
V25kD2L8exp	603.90	212.05	616.39	210.77	307.41	211.70	642.46	1,384.24
V25kD2L8norm	370.65	953.96	389.75	1,020.88	168.51	942.32	405.34	1,151.69
V25kD2L8uni	315.86	897.93	311.52	950.25	160.47	1,003.63	335.32	1,823.13
V125kD2L8exp	381.12	3,119.93	386.20	3,233.09	231.98	3,282.21	475.99	5,386.13
V125kD2L8norm	334.58	4,533.95	340.76	4,828.57	221.42	4,793.20	342.29	5,515.67
V125kD2L8uni	189.58	3,179.83	187.02	3,308.09	125.31	3,181.34	213.12	3,863.40
advogato	1,032.94	5.01	1,027.80	3.69	625.76	3.69	1,207.44	21.32
yagoFacts-small	8.75	10.67	9.11	11.40	4.44	12.67	9.00	10.05
subeljCoraL8exp	788.96	142.83	812.33	140.80	379.78	141.93	947.71	409.45
arXivhepphL8exp	865.54	4.52	902.63	4.31	563.99	4.30	1,597.88	15.14
p2p-GnutellaL8exp	759.17	4.09	794.94	4.36	439.87	4.37	846.21	770.44
Average	510.51	784.25	522.49	1,035.80	525.91	771.16	591.83	1,401.42

Table 6.14: Total speed-ups for all 6 query conditions and PA-datasets using LI+OTH+EXTv1 ($k = n/20, b = 20$).

degree, $ \mathcal{L} $	qs_0, true	qs_0, false	qs_1, true	qs_1, false	qs_2, true	qs_2, false
	$ \mathcal{L} /4$		$ \mathcal{L} /2$		$ \mathcal{L} - 2$	
3, 8	575.65	34.72	491.32	295.45	565.15	711.61
3, 10	629.00	2.38	540.14	307.73	571.57	816.15
3, 12	666.71	21.72	724.43	393.65	585.70	957.84
3, 14	710.23	6.94	682.77	254.04	630.65	846.47
3, 16	977.86	3.00	535.27	10.50	780.25	40.88
4, 8	550.42	113.30	628.94	8.30	649.86	129.08
4, 10	590.96	16.92	508.76	50.65	647.62	485.68
4, 12	595.97	3.82	605.84	514.64	445.91	119.14
4, 14	552.87	62.47	497.39	518.51	493.54	28.38
4, 16	-	-	-	-	-	-
5, 8	778.94	297.72	614.06	699.38	795.42	64.36
5, 10	682.26	147.58	417.01	827.39	682.71	16.51
5, 12	550.57	5.20	525.59	772.23	436.84	14.68
5, 14	664.51	3.22	790.71	640.55	613.85	11.40
5, 16	-	-	-	-	-	-

Table 6.15: The datasets included in this part of the experiment with k expressed as the fraction of the number of vertices n .

Dataset	k	n
jmd	$n/50$	486,320
citeSeerL8exp	$n/50$	384,414
NotreDameL8exp	$n/50$	325,730
soc-sign-epinionsL8exp	$n/50$	131,828
socSlashdotL8exp	$n/50$	82,168
TwitterL8exp	$n/100$	465,018
webBerkstanL8exp	$n/100$	685,231
webStanfordL8exp	$n/100$	281,904
V125kD5L8exp	$n/100$	125,000
pl125ka2.0L8exp	$n/100$	125,000
webGoogleL8exp	$n/250$	875,713
zhishihudongL8exp	$n/250$	2,452,715
usPatentsL8exp	$n/250$	3,774,769
V625kD5L8exp	$n/250$	625,000
pl625ka2.0L8exp	$n/250$	625,000
socPokec	$n/500$	1,632,803

6.4 Part 3: large graphs ($|E| > 500,000$)

Table 6.15 shows the datasets divided into four categories where k is either $\frac{n}{50}$, $\frac{n}{100}$, $\frac{n}{250}$ or $\frac{n}{500}$. These distinctions were made to not let exceed any of the datasets the 6 hours time limit and the 128GB memory limit while still being able to provide a reasonable speed-up.

For $n/50$ we used LI+OTH with $b = 20$ and LI+OTH+EXTv1 with $b = 20$ and $b = 40$ to demonstrate the effect of the first extension here. For the remainder we only used LI+OTH+EXTv1 with $b = 20$ and $b = 40$. We do this to find out whether a larger budget works better or not.

6.4.1 LI+OTH+EXTv1: Index construction time (s) and size (MB)

In Table 6.16 we see the index construction time (s) and size (MB) for several datasets and LI+OTH+EXTv1 $b = 20$ and $b = 40$. We excluded the data for LI+OTH from the top section of the table as it is very similar to LI+OTH+EXTv1.

There is very little to no difference between $b = 20$ and $b = 40$. This surprising as any non-landmarked vertex $v_0 \in V$ needs to do a double amount of work in the second case. However when $b = 40$ any $v \in V$ (landmark or not) that has been indexed, can be used in forward propagation. Hence with $b = 40$ there are more opportunities for using FORWARDPROP, which might explain why there is hardly a difference in the index construction time between the two.

socPokec did not complete because it exceeded the maximum memory size of 128GB. By then it had built around 70% of the landmarks within around 16,000 seconds, which is at roughly 75% of the maximal time.

jmd has a very small index, because it is acyclic. The same holds for **usPatentsL8exp**. A large chunk of the index size for both is due to the graph size as well. It is difficult to generate difficult queries for these kinds of datasets. A difficult query is a query for which BFS needs to visit a sufficiently large part of the vertices. As we choose a random starting point for each query, we might end up at a point where visiting enough vertices is not possible.

The datasets in the upper section of Table 6.16 have their indices built within 800 seconds each. This is in sharp contrast with the datasets in the second and third section of the table. The datasets in the first section have a very skewed out-degree distribution, which favours LI. The datasets in the third section have a large number of edges as well.

Table 6.16: Index construction time (s) and index size (MB) for all the datasets involved.

Dataset	Time (s)	Size (MB)	Time (s)	Size (MB)
	LI+OTH+EXTv1 $b = 20$		LI+OTH+EXTv1 $b = 40$	
jmd	468.91	22.55	538.28	21.74
NotreDameL8exp	743.47	7,390.52	740.82	7,387.86
citeSeerL8exp	179.4	2,120.12	180.25	2,154.19
soc-sign-epinionsL8exp	359.05	3,426.55	358.73	3,425.45
socSlashdotL8exp	258.01	3,686.17	257.68	3,684.49
plV125ka2.0L8exp	158.31	1,988.32	157.83	1,987.27
V125kD5L8exp	1,483.89	12,395.29	1,483.46	12,394.51
TwitterL8exp	1,526.82	17,597.71	1,527.14	17,588.68
webBerkstanL8exp	8,555.18	54,040.7	8,549.91	54,034.57
webStanfordL8exp	3,929.01	24,818.4	3,926.55	24,816.83
zhishihudongL8exp	9,971.12	29,738.07	9,958.60	29,736.57
webGoogleL8exp	8,351.06	53,923.01	8,347.77	53,919.79
V625kD5L8exp	12,400.35	66,290.43	12,404.22	66,644.64
plV625ka2.0L8exp	1,955.14	19,446.65	1,960.67	19,441.35
usPatentsL8exp	7,539.27	705.02	7,503.74	563.98
socPokec	-	-	-	-
Average	3,858.6	19,839.3	3,859.71	19,853.46

In Figures 6.4 and 6.5 we see the development of the index construction time and the index size as a function of the percentage of landmarks built. This was only done for the datasets **plV625ka2.0L8exp** and **webGoogleL8exp**. The two datasets exhibit a different development for the index construction time. **plV625ka2.0L8exp** has a more linear development, whereas **webGoogleL8exp** has a stronger development in the beginning. W.r.t. the index size both datasets exhibit a similar development.

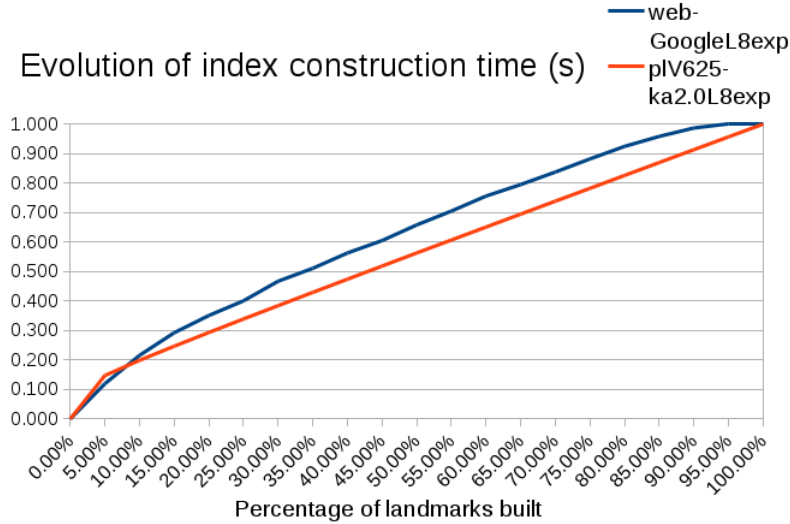


Figure 6.4: On the y-axis we see the percentage of the total time needed to build all landmarks against the percentage of landmarks (LI+OTH+EXTv1) that have been built on the x-axis for datasets **plV625ka2.0L8exp** and **webGoogleL8exp**.

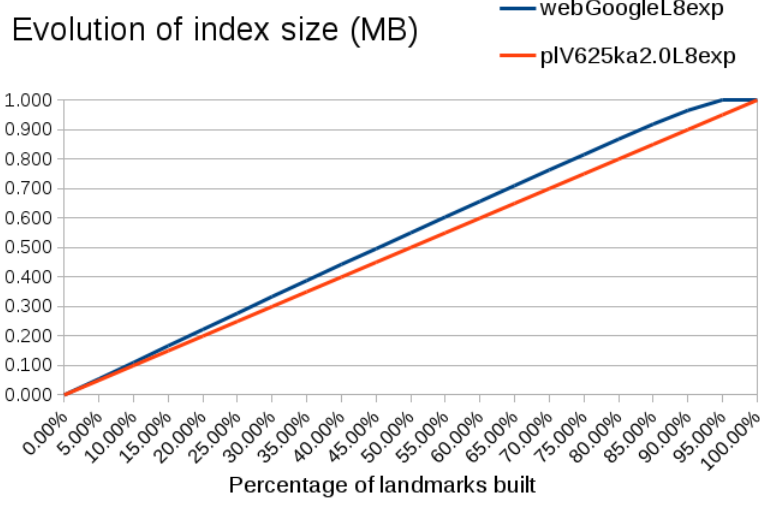


Figure 6.5: On the y-axis we see the percentage of the total memory needed to build all landmarks against the percentage of landmarks (LI+OTH+EXTv1) that have been built on the x-axis for datasets **plV625ka2.0L8exp** and **webGoogleL8exp**.

d

6.4.2 LI+OTH+EXTv1: Total speed-ups

Table 6.17: Mean query execution time (ms) of BFS for all datasets per query condition.

dataset	qs_1, true $ \mathcal{L} /4$	qs_2, true $ \mathcal{L} /2$	qs_3, true $ \mathcal{L} - 2$	qs_1, false	qs_2, false	qs_3, false
jmd	0.003	0.003	0.003	0.003	0.003	0.003
NotreDameL8exp	0.237	0.183	0.194	0.931	1.826	1.070
soc-sign-epinionsL8exp	3.818	5.166	4.786	6.930	8.610	9.846
socSlashdotL8exp	3.355	5.345	3.547	7.666	10.641	9.047
citeSeerL8exp	0.861	1.030	1.254	1.582	2.122	2.619
plV125ka2.0L8exp	4.909	4.780	5.171	4.720	5.773	7.466
TwitterL8exp	5.608	5.049	5.826	7.873	7.166	1.785
V125kD5L8exp	0.281	0.165	0.221	6.764	5.380	5.745
webBerkstanL8exp	5.649	5.051	3.464	48.107	30.779	48.158
webStanfordL8exp	4.723	3.891	4.694	16.000	12.346	20.791
V625kD5L8exp	1.065	1.240	1.424	7.304	34.741	75.612
plV625ka2.0L8exp	45.528	47.592	47.831	52.401	63.567	70.526
usPatentsL8exp	0.324	0.407	0.422	0.332	0.416	0.391
webGoogleL8exp	31.727	26.746	26.052	90.562	50.733	89.902
zhishihudongL8exp	10.149	13.326	11.869	36.621	34.458	34.125
Average	7.882	7.998	7.784	19.186	17.904	25.139

Table 6.17 shows the mean query answering times of BFS. Table 6.18 shows the total speed-ups for the datasets for which we set $k = \frac{n}{50}$ and the methods LI+OTH+EXTv1 with a budget $b = 20$ (first row), LI+OTH+EXTv1 with $b = 40$ (second row) and LI+OTH with $b = 20$ (third row).

Table 6.18: Speed-ups for the methods LI+OTH+EXTv1 with a budget $b = 20$ (first row), LI+OTH+EXTv1 with $b = 40$ (second row) and LI+OTH with $b = 20$ (third row) and the datasets for which we set $k = \frac{n}{50}$.

dataset	$qs_{1,true}$ $ \mathcal{L} /4$	$qs_{1,false}$	$qs_{2,true}$ $ \mathcal{L} /2$	$qs_{2,false}$	$qs_{3,true}$ $ \mathcal{L} - 2$	$qs_{3,false}$
jmd	7.18	1.27	2.63	1.21	1.16	0.98
	6.58	1.08	2.39	1.08	1.14	0.92
	7.22	1.28	2.63	1.22	1.18	0.98
NotreDameL8exp	142.24	43.92	107.23	31.07	113.35	16.86
	135.8	17.54	105.9	8.61	110.29	12.81
	143.42	17.81	108.08	8.74	113.29	13.03
citeSeerL8exp	382.66	123.29	440.74	98.4	552.45	84.3
	254.44	86.68	394.09	94.84	571.69	81.23
	386.32	123.75	443.33	98.34	563.44	84.02
soc-sign-epinionsL8exp	5,372.24	4.77	6,973.25	4.36	6,801.43	3.4
	5,149.79	4.48	6,417.96	2.94	6,132.68	2.9
	5,364.51	4.57	7,019.03	3.01	6,780.85	2.99
socSlashdotL8exp	5,302.07	2.42	9,942.41	2.46	6,666.50	3.98
	4,215.56	2.09	9,016.06	1.95	6,137.54	2.58
	5,289.85	2.15	10,110.35	2.01	6,647.29	2.67
Average	2,241.28	35.13	3,493.25	27.5	2,826.98	21.9
	1,952.43	22.37	3,187.28	21.88	2,590.67	20.09
	2,238.26	29.91	3,536.68	22.66	2,821.21	20.74

jmd is a special case. It is acyclic, like **usPatentsL8exp**. This makes the mean query answering times much lower, as it is more difficult to generate difficult queries for acyclic datasets. Because BFS performs relatively well on acyclic datasets, we get low total speed-ups for **jmd**.

NotreDameL8exp has significant speed-ups for all query conditions. There is a clear difference between the first row (LI+OTH+EXTv1 with $b = 20$) and the third row here (LI+OTH with $b = 20$) for **Notre-DameL8exp**. The first extension pays off here clearly. However adding more budget to each non-landmark node, i.e. comparing the first and the second row, has a negative effect. Similar results can be seen for **soc-sign-epinionsL8exp** and **socSlashdotL8exp**.

Table 6.19 shows the speed-ups for the datasets for which we had set $k = \frac{n}{100}$. The speed-ups for **TwitterL8exp** are very high, compared to the other speed-ups and even the speed-ups before. The reason for this is that **TwitterL8exp** has an extremely skewed out-degree distribution and that we select the landmarks based on their total degree (in- plus out-degree). Hence any query (v, w, L) will hit a landmark v' within a few steps. As the landmarks have a lot of out-edges as well, a relative large space is pruned whenever direct attempt (v', w, L) is false. Also note that the mean BFS query times for **TwitterL8exp** are also on the high side looking at Table 6.16. **TwitterL8exp** is an exception though if we look at the speed-ups.

Table 6.20 shows the speed-ups for the datasets for which we had set $k = \frac{n}{250}$. The results for the True-queries are still very good, but for the False-queries the results are bad, often close to no total speed-up at all.

The asymmetry between the speed-ups for the True- and False-queries is growing when we compare Table 6.18 with Table 6.19 and Table 6.19 with Table 6.20. This is because the ratio $\frac{k}{n}$ is decreasing. As the graph and the query difficulty grow, one can imagine that resolving a query (v, w, L) directly through a landmark $v' \in V'$ has more of an effect. True-queries have higher total speed-ups, because in practice there will often be multiple landmarks $v' \in V'$ that have a L -path to w and hitting one of these is sufficient. False-queries do not have this advantage. Moreover, the first extension of LI only works for the first direct attempt. In the experiments leading up to this experiment we tried doing this for multiple landmarks, but this did not yield any better results. A

Table 6.19: Speed-ups for the methods LI+OTH+EXTv1 with a budget $b = 20$ (first row) and LI+OTH+EXTv1 with $b = 40$ (second row) and the datasets for which we set $k = \frac{n}{100}$.

dataset	qs_1, true $ \mathcal{L} /4$	qs_1, false	qs_2, true $ \mathcal{L} /2$	qs_2, false	qs_3, true $ \mathcal{L} - 2$	qs_3, false
plV125ka2.0L8exp	6,853.96 6,382.26	6.43 6.29	7,296.94 7,260.51	7.94 7.85	8,102.48 7,818.16	9.42 9.26
TwitterL8exp	56,020.11 56,079.89	27,872.16 22,972.38	50,487.78 50,326.67	71,657.92 71,657.92	58,264.06 58,264.06	17,852.12 17,852.12
V125kD5L8exp	336.24 263.19	1.22 1.17	244.2 222.46	1.91 1.29	291.17 256.7	10.55 9.34
webBerkstanL8exp	1,437.66 1,267.31	5.98 5.83	1,217.78 1,135.14	4.41 4.3	905.13 870	4.19 4.07
webStanfordL8exp	3,299.61 3,266.36	25.45 21.83	2,566.25 2,700.26	22.88 19.45	3,099.08 3,188.28	14.56 12.02
Average	13,589.52 13,451.8	5,582.25 4,601.5	12,362.59 12,329.01	14,339.01 14,338.16	14,132.38 14,079.44	3,578.17 3,577.36

better version of the first extension is necessary to make the asymmetry between the two smaller.

Table 6.20: Speed-ups for the methods LI+OTH+EXTv1 with a budget $b = 20$ (first row) and LI+OTH+EXTv1 with $b = 40$ (second row) and the datasets for which we set $k = \frac{n}{250}$.

dataset	qs_1, true $ \mathcal{L} /4$	qs_1, false	qs_2, true $ \mathcal{L} /2$	qs_2, false	qs_3, true $ \mathcal{L} - 2$	qs_3, false
zhishihudongL8exp	775.96 749.67	1 1	999.37 1,026.45	0.89 0.9	891.01 884.79	0.99 1
V625kD5L8exp	320.04 322.75	1,572.78 1,611.81	361.58 371.02	3.07 3.05	402.29 408.75	1.22 1.21
plV625ka2.0L8exp	14,862.82 14,808.65	5.09 5.09	15,426.74 15,512.66	5.25 5.24	15,565.15 15,475.5	2.07 2.07
usPatentsL8exp	1.73 1.73	1.31 1.31	3.03 3.05	1.38 1.38	1.56 1.59	1.1 1.1
webGoogleL8exp	6,379.07 6,442.69	1.42 1.4	3,970.09 3,741.13	2.48 2.44	5,904.38 5724.36	1.98 1.95
Average	4,467.92 4,465.1	316.32 324.12	4,152.16 4,130.86	2.61 2.6	4,552.88 4,499	1.47 1.47

6.4.3 LI+OTH+EXTv2: Index construction time (s) and size (MB)

The disappointing results for the False-queries in Table 6.20 lead us to re-examine the first extension, as this extension was particularly aimed at speeding up the False-queries. We found a different way of using the idea behind the first extension which ended up being called the third extension, i.e. LI+OTH+EXTv2.

Table 6.21 shows the same datasets as Table 6.16 but with a different often lower number of landmarks. The experiments were conducted with the methods LI+OTH+EXTv2 with $b = 15$ and $b = 30$ and LI+OTH+EXTv1 with $b = 15$.

Looking at Table 6.22 we see there is a small difference between the LI+OTH+EXTv2 and LI+OTH+EXTv1 methods when we look at the index size and index construction time. Often there is only a few hundreds of MB

Table 6.21: The datasets included in this part of the experiment with k expressed as the fraction of the number of vertices n .

Dataset	k	n
jmd	$n/100$	486,320
citeSeerL8exp	$n/100$	384,414
NotreDameL8exp	$n/100$	325,730
soc-sign-epinionsL8exp	$n/100$	131,828
socSlashdotL8exp	$n/100$	82,168
V125kD5L8exp	$n/100$	125,000
pl125ka2.0L8exp	$n/100$	125,000
TwitterL8exp	$n/100$	465,018
webBerkstanL8exp	$n/500$	685,231
webStanfordL8exp	$n/500$	281,904
webGoogleL8exp	$n/500$	875,713
zhishihudongL8exp	$n/500$	2,452,715
usPatentsL8exp	$n/500$	3,774,769
V625kD5L8exp	$n/500$	625,000
pl625ka2.0L8exp	$n/500$	625,000

Table 6.22: Index construction time (s) and size (MB) for LI+OTH+EXTv2 with $b = 15$ or $b = 30$ and LI+OTH+EXTv1 with $b = 15$.

dataset	LI+OTH+EXTv2 $b = 15$		LI+OTH+EXTv2 $b = 30$		LI+OTH+EXTv1 $b = 15$	
citeSeerL8exp	114.03	1,158.75	115.15	1,181.46	97.63	1,000.99
jmd	23.72	21	23.47	21.05	3.92	18.4
NotreDameL8exp	395.12	4,777.49	394.96	4,782.76	344.13	4,112.74
plV125kL8a2.0L8exp	183.99	2,538.03	184.7	2,538.97	158.85	1,987.15
V125kD5L8exp	957.97	6,467.02	958.84	6,523.72	950.78	6,353.11
soc-sign-epinionsL8exp	205.51	2,390.25	205.9	2,395.20	171.64	1,713.02
socSlashdotL8exp	138.25	2,137.47	138.64	2,146.08	126	1,853.14
TwitterL8exp	889.16	17,632.80	889.17	17,632.80	866.23	17,427.31
zhishihudongL8exp	6,418.79	16,198.93	6,435.23	16,572.12	5,623.60	14,802.09
webBerkstanL8exp	1,604.28	9,879.21	1,603.63	9,888.58	1,580.64	9,476.28
webGoogleL8exp	4,690.88	27,117.16	4,694.70	27,147.40	4,736.21	26,813.72
webStanfordL8exp	295.57	2,557.07	295.98	2,557.07	292.12	2,499.10
usPatentsL8exp	6,463.42	21,915.31	6,399.90	21,915.31	5,561.24	475.39
plV625kL8a2.0L8exp	1,233.28	12,982.43	1,240.02	12,982.43	1,063.14	9,726.56
V625kD5L8exp	6,433.90	34,300.96	6,437.99	34,568.91	6,897.07	32,534.70
Average	2,003.19	10,804.93	2,001.22	10,856.92	1,898.21	8,719.58

difference between LI+OTH+EXTv2 and LI+OTH+EXTv1, which is relatively low. The difference in the construction time is often larger but manageable, i.e. no more than 15% increase. Only for **jmd** the difference is quite large, i.e. 7 times larger.

All datasets have an index construction time below (6,500 seconds). The datasets in the upper section of the table all have construction times below 1,000 seconds. The index size relative to the number of vertices is also low. Considering the space complexity for k landmarks, i.e. $O(b + |V| \cdot k \cdot 2^{|\mathcal{L}|})$ (see Section 4.2.2) we have a very small index size.

6.4.4 LI+OTH+EXTv2: Total speed-ups

Table 6.23 shows the speed-ups for the datasets with $k = \frac{n}{100}$. Note that the query times of BFS are still the same (see Table 6.17). In the first two rows of each dataset we see LI+OTH+EXTv2 with $b = 15$ and $b = 30$ and in the third row we see LI+OTH+EXTv1 with $b = 15$. The speed-ups for the True-queries are roughly comparable for the first two rows and the third row, but for the False-queries there is a clear improvement. The magnitude of the improvement differs per dataset.

When we compare the $b = 15$ -version against the $b = 30$ -version there is no clear winner here. In some occasions one variant works better and in some occasions the other works better.

Table 6.23: Speed-ups for LI+OTH+EXTv2 $b = 15$ or $b = 30$ and LI+OTH+EXTv1 $b = 15$ and the datasets for which we set $k = \frac{n}{100}$.

dataset	$qs_{1,true}$ $ \mathcal{L} /4$	$qs_{1,false}$	$qs_{2,true}$ $ \mathcal{L} /2$	$qs_{2,false}$	$qs_{3,true}$ $ \mathcal{L} - 2$	$qs_{3,false}$
citeSeerL8exp	208.53	82.34	279.52	76.05	389.55	55.04
	229.93	81.37	308.25	73.83	443.08	54.91
	211.47	2.23	294.43	1.71	407.37	1.33
jmd	7.22	0.98	1.75	0.95	1.1	1
	6.78	0.93	1.72	0.96	1.11	1.01
	7.14	0.99	1.74	0.98	1.12	0.99
NotreDameL8exp	14.55	91.55	16.51	102.99	65.22	215.55
	14.49	87.92	16.29	165.59	65.66	235.31
	14.54	4.02	16.24	2.26	65.56	1.77
plV125kL8a2.0L8exp	1,498.19	418.39	1,875.59	671.07	2,561.84	445.03
	1,655.61	415.73	1,973.44	689.51	3,101.36	472.22
	1,489.66	1.7	1,915.09	2.17	2,559.05	2.76
soc-sign-epinionsL8exp	1,733.57	1,894.27	3,930.49	2,755.08	4,213.19	2,958.56
	1,769.20	1,926.03	4,128.76	2,758.99	4,117.01	2,927.49
	1,736.41	1.52	3,893.73	1.27	4,145.18	1.49
socSlashdotL8exp	1,065.00	1,274.13	5,971.36	1,992.09	1,503.16	656.8
	1,049.51	1,258.57	5,794.21	1,936.08	1,514.89	650.81
	1,051.95	1.24	5,961.80	1.26	1,476.11	1.37
TwitterL8exp	54,335.96	27,718.68	48,976.11	69,439.02	56,912.68	17,487.74
	54,335.96	27,289.48	48,976.11	69,439.02	56,912.68	17,487.74
	54,335.96	27,211.45	48,976.11	69,439.02	56,912.68	17,487.74
Average	8,409.0	4,497.19	8,721.62	10,719.61	9,378.11	3,117.1
	8,437.35	4,437.15	8,742.68	10,723.43	9,450.83	3,118.5
	8,406.73	3,889.02	8,722.73	9,921.24	9,366.72	2,499.64

Table 6.24 shows the total speed-ups but then for $k = \frac{N}{500}$. Again we see a strong performance improvement for LI+OTH+EXTv2 with respect to the False-queries. The improvement is the highest for $|\mathcal{L}|/4$ which has to do with the following. When we prune a subset of the graph by using the third extension, the effect of this is larger of queries with a smaller label set because a smaller number of the edges in the graph can be used. A way to get higher numbers for the False-queries in case $|L| = |\mathcal{L}| - 2$ would be to increase the *MAXDIST*-parameter, currently equal to $|sL|/4 + 1$, of the third extension. This would include more entries into *seqE* and give queries with a larger label set the opportunity to prune even more vertices.

The total speed-ups for **usPatentsL8exp** are bad. We saw the same thing with **jmd**. This is caused by the difficulty to generate difficult queries for acyclic graphs, which makes BFS perform relatively well. Another part may have to do with the way we build the index. For acyclic graphs we would prefer a different construction algorithm. For instance, build a landmark for every root of a subtree T' where the number of vertices in T' is

Table 6.24: Speed-ups for LI+OTH+EXTv2 $b = 15$ or $b = 30$ and LI+OTH+EXTv1 $b = 15$ and the datasets for which we set $k = \frac{n}{500}$.

dataset	$qs_{1,true}$ $ \mathcal{L} /4$	$qs_{1,false}$	$qs_{2,true}$ $ \mathcal{L} /2$	$qs_{2,false}$	$qs_{3,true}$ $ \mathcal{L} - 2$	$qs_{3,false}$
zhishihudongL8exp	803.99	911.96	1,057.82	106.22	954.23	20.77
	804.83	911.62	1,062.50	106.35	954.31	20.71
	803.63	0.93	1,056.27	0.98	953.79	0.93
webBerkstanL8exp	517.56	342.29	228.37	31.49	894.44	40.62
	513.55	341.16	278.72	31.63	877.37	40.69
	518.5	1.21	228.94	1.4	895.1	1.32
webGoogleL8exp	4,181.61	5,908.13	2,246.65	17.12	4,385.16	20
	4,310.40	6,060.13	2,447.64	17.59	4,638.08	20.24
	4,146.34	1.08	2,252.71	1.22	4,388.31	1.25
webStanfordL8exp	1,414.59	239.56	295.47	15.82	2,380.15	13.06
	1,381.55	241.09	820.96	15.85	2,446.68	14.4
	1,408.43	1.24	293.87	1.13	2,354.91	1.19
usPatentsL8exp	1.12	1.04	1.4	1.16	1.11	0.98
	1.11	1.02	1.42	1.15	1.13	0.97
	1.09	1	1.34	1.11	1.09	0.98
plV625kL8a2.0L8exp	3,913.11	1,257.15	2,416.38	709.51	2,304.48	246.42
	3,924.00	1,233.11	2,981.31	700.4	2,303.75	245.42
	3,921.50	1.9	2,417.68	1.63	2,305.84	1.16
V625kD5L8exp	183.49	56.56	218.79	10.93	270.66	3.09
	188.2	58.1	254.12	11.24	272.59	3.26
	183.64	1.59	217.75	0.72	268.83	0.67
Average	1,573.64	1,245.24	923.55	127.46	1,598.6	49.28
	1,589.09	1,263.75	1,120.95	126.32	1,641.99	49.38
	1,569.02	1.28	924.08	1.17	1,595.41	1.07

larger than $\frac{n}{K}$ for some constant K .

6.5 Maintenance

We tested the correctness and speed of our maintenance methods by doing the following. First we built a LI I_1 for a graph G . Then we applied $K = 30$ random operations (adding an edge, removing an edge) to G and to I_1 . Afterwards, we built another LI I_2 . To test the correctness of our approach we compared I_1 against I_2 over $L = 200$ random queries. To assess the speed of our approach we calculated the time needed for any of the K operations as a ratio of the time to build I_1 . This ratio should be a low value.

6.5.1 Adding an edge

Table 6.25 shows for 4 datasets the percentage of the construction time of I_1 (before the updates) that it took to add a random edge (v, w, l) to the graph and update I_1 . For example, “Average 0.17” means that K insertions of a new edge took on average 17% of the initial index construction time.

We can see that on average for label sets with an exponential distribution inserting a new edge takes relatively more time (roughly 22%) than for datasets with a normal or uniform distribution (4% and 6%). The data in the table shows that frequent insertion of edges on a large index is not desirable. A more efficient method is thus desired.

Table 6.25: To each of the graphs listed in the columns we added 30 random new edges. The ratio of the time needed to update the index for any new edge and the original index construction time, i.e. before adding any of the 30 edges, was calculated. The average and standard deviation of these 30 ratios are reported in this table.

	V5kD2L8uni	V5kD2L8norm	V5kD2L8exp	V25kD2L8exp
Average	0.04	0.06	0.17	0.22
Standard deviation	0.03	0.03	0.11	0.11

6.6 Extensions

6.6.1 Query for all nodes

A QueryAll-query (v, L) returns a subset of vertices $V' \subseteq V$ where $v' \in V$ if and only if query (v, v', L) is true. We ran 200 queries (half of which $|L| = |\mathcal{L}|/2$ and half of which $|L| = |\mathcal{L}| - 2$) for several datasets. We selected queries such that each query would hit at least 10% of the vertices. For uni- or norm-datasets this requirement would not be met if we had set $|L| = |\mathcal{L}|/4$. Hence we omitted these types of queries.

Figures 6.6 and 6.7 show the individual speed-up for QueryAll-queries in two different settings. We varied the label set distribution in the first setting. In the second setting we varied the degree. From the first figure we can see that the individual speed-ups are higher for an exp-dataset than for a uni- or norm-dataset. The degree is not that influential.

Table 6.26: Average speed-up for QueryAllQueries.

dataset	speed-up $ L = \mathcal{L} /2$	speed-up $ L = \mathcal{L} - 2$
V5kD2L8uni	1.24	2.49
V5kD2L8norm	1.40	2.90
V5kD2L8exp	5.97	8.86
V25kD3L8exp	4.40	7.55
V25kD4L8exp	4.04	6.95
V25kD5L8exp	3.81	6.80

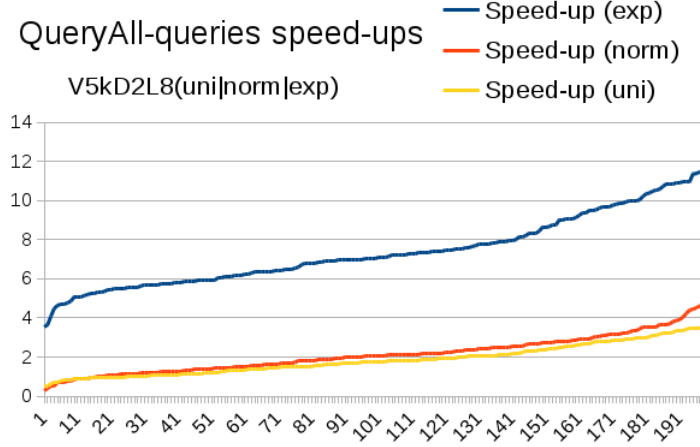


Figure 6.6: Speed-ups per query for QueryAll-queries for PA -datasets with 5,000 vertices, a degree of 2, 8 labels and either a uniform, normal or exponential label set distribution. We used LI+OTH+EXTv1($k = N/10, b = 0$).

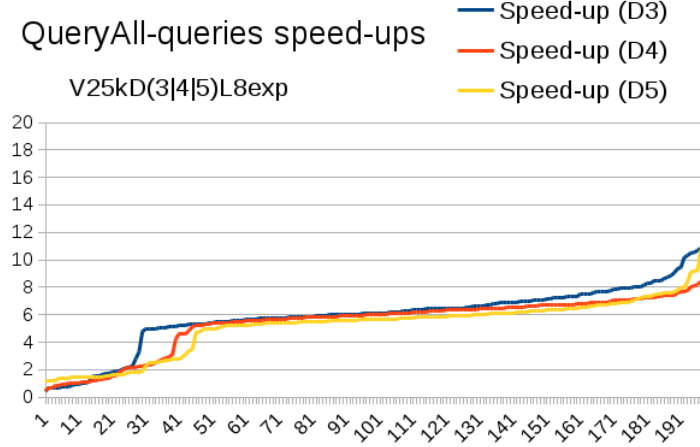


Figure 6.7: Speed-ups per query for QueryAll-queries for PA -datasets with 25,000 vertices, a degree of 3, 4 or 5, 8 labels and an exponential label set distribution. We used LI+OTH+EXTv1($k = N/10, b = 0$).

6.6.2 Distance queries

Table 6.27 shows the index size and construction time for two versions of LI+OTH+EXTv1 ($k = n/10, b = 0$) and some datasets. The first version is the normal “LCR”-version. The second version is the version that can give the distance of the shortest path P for a query (v, w, L) s.t. $\text{Labels}(P) \subseteq L$.

The results show that the index size is always at least 2 times larger. The index construction time can even grow much stronger. We do wish to note that this solution returns the exact distance, i.e. it does not return an approximation of the distance.

Table 6.27: Index size (MB) and index construction time (s) for LI+OTH+EXTv1 and LI+OTH+EXTv1 with distance ($k = N/10, b = 0$).

dataset	size (MB)	time (s)	size (MB) WD	time (s) WD	ratio (size)	ratio (time)
V5kD2L8uni	84.7	8.4	380.6	65.4	4.5	7.7
V5kD2L8norm	71.6	5.4	307.7	39.4	4.3	7.3
V5kD2L8exp	35.2	102.3	0.77	4.5	2.9	5.9
V25kD2L8exp	857.1	20.4	2,871.9	217.6	3.3	10.6
advogato	54.5	1.7	110.1	58.1	2.0	33.9

Chapter 7

Conclusion

In this paper we studied the reachability problem with a label-constraint. This problem has a lot of potential in real world applications (RDF, social networks and biological networks).

Several methods (LI, PARTIALINDEX, DOUBLEBFS, NEIGHBOUREXCHANGE, CLUSTEREDEXACT and ZOU) were experimentally compared against each other. The results clearly were in favour of the new approach we have introduced in this thesis (LI). This method proved to have a significant improvement over BFS considering query answering times. In combination with the second and third extension, i.e. LI+OTH+EXTv2, we were able to achieve excellent speed-ups for large datasets using a relatively low number of landmarks. However datasets with a close to uniform out-degree distribution (like ER-graphs) or datasets with a large label set ($|\mathcal{L}| \geq 16$) still remain difficult. Using just BFS in these cases might be a good solution as this method has an acceptable query answering time.

We can draw the following conclusions in the end.

1. Our major contribution LI+OTH+EXTv2 achieves excellent speed-ups. There is an asymmetry between true- and false-queries. True-queries can stop after finding their target. False-queries have to explore larger sections of the graph.
2. When the ratio between the number of landmarks k and the number of vertices $n \frac{k}{n}$ is low, then the first and the third extension prove their value and we can still achieve excellent speed-ups. Hence the approach LI+OTH+EXTv2 is scalable.
3. The (out)-degree distribution is an important parameter for the index construction time and size of LI+OTH+EXTv2.
4. The label set size $|\mathcal{L}|$ of the labelled graph and the degree per node $D = \frac{n}{m}$ are also two important parameters for the index construction time and size of LI+OTH+EXTv2. There is an interaction effect between the two parameters.
5. Index maintenance for LI+OTH+EXTv2 remains an issue. It is difficult to estimate which entries in an index are no longer valid after an update and to create an algorithm that is significantly faster than rebuilding index.
6. Alternative approaches (PARTIALINDEX, DOUBLEBFS, NEIGHBOUREXCHANGE, CLUSTEREDEXACT and ZOU) are not competitive to LI+OTH+EXTv2 w.r.t. their speed-ups and/or their index construction times.

7.1 Future work

One possible extension could be to use our main contribution, i.e. LI+OTH+EXTv2, to also answer richer queries. One of these could be distance queries. We have only built an index that can answer these types of queries, but we have not looked into the query algorithm for this. A solution to this problem could return an exact distance or an approximation of that distance.

Another query could be to find a witness of a query (v, w, L) , i.e. an actual path P that connects v and w .

A different way to pick the landmarks is also still an open issue. One could choose to use some centrality measure instead of the total degree.

The index construction time is still relatively high for the large datasets. It might not be possible to reduce it further without paying some price in terms of memory or speed-up.

A good maintenance algorithm is still on the list. Addition of edges is still very costly. Removal of edges might require adding some redundancy, e.g. storing some pairs $(u, L'), (u, L) \in \text{Ind}(v)$ for $v \in V$ s.t. $L' \subseteq L$.

One could still look at the potential of making a multi-threaded version of LI. We did not have the time to study this.

Bibliography

- [1] Francesco Bonchi, Aristides Gionis, Francesco Gullo, and Antti Ukkonen. Distance oracles in edge-labeled graphs. In *Proceedings of the International Conference on Extending Database Technology (EDBT)*, pages 547–548, Athens, 2014. 2, 3, 9, 35
- [2] Minghan Chen, Yu Gu, Yubin Bao, and Ge Yu. *Label and Distance-Constraint Reachability Queries in Uncertain Graphs*. Lecture Notes in Computer Science. Springer International Publishing, 2014. 5
- [3] Edith Cohen, Eran Halperin, Haim Kaplan, and Uri Zwick. Reachability and distance queries via 2-hop labels. *SIAM J Comput*, 32(5):1338–1355, 2003. 7
- [4] George Fletcher and Yuichi Yoshida. Notes on landmark labeling label-constrained reachability queries. Unpublished, 2015. 8, 9, 16, 17
- [5] Ruoming Jin, Hui Hong, Haixun Wang, Ning Ruan, and Yang Xiang. Computing label-constraint reachability in graph databases. In *ACM SIGMOD*, pages 123–134, Indianapolis, 2010. 5
- [6] Jrme Kunegis. Konect - the koblenz network collection. In *Proc. Int. Conf. on World Wide Web Companion*, pages 1343–1350, Koblenz, 2013. 37
- [7] Jure Leskovec and Andrej Krevl. SNAP Datasets: Stanford large network dataset collection. <http://snap.stanford.edu/data>, June 2014. 37
- [8] Jure Leskovec and Rok Sosič. SNAP: A general purpose network analysis and graph mining library in C++. <http://snap.stanford.edu/snap>, June 2014. 37
- [9] Yosuke Yano, Takuya Akiba, Yoichi Iwata, and Yuichi Yoshida. Fast and scalable reachability queries on graphs by pruned labeling with landmarks and paths. In *ACM CIKM*, pages 1601–1606, San Francisco, 2013. 3
- [10] Jeffrey Xu Yu and Jiefeng Cheng. *Graph Reachability Queries: A Survey*. Advances in Database Systems. Springer US, 2010. 3, 7, 8
- [11] Lei Zou, Kun Xu, Jeffrey Xu Yu, Lei Chen, Yanghua Xiao, and Dongyan Zhao. Efficient processing of label-constraint reachability queries in large graphs. *Inf. Syst.*, 40:47–66, 2014. 2, 5, 9, 11, 14, 15, 20, 23