

“

COMPACT REPRESENTATION OF SEPARABLE GRAPH

by

Xiang Zhang

Submitted in partial fulfillment of the requirements
for the degree of Master of Applied Computer Science

at

Dalhousie University
Halifax, Nova Scotia
August 2017

© Copyright by Xiang Zhang, 2017

Contents

| | |
|---|------------|
| List of Tables | iii |
| List of Figures | iv |
| Abstract | v |
| Acknowledgements | vi |
| Chapter 1 Introduction | 1 |
| 1.1 Related Works | 4 |
| 1.2 Preliminaries | 6 |
| 1.3 Outline | 10 |
| Chapter 2 Edge Separator Tree | 11 |
| Chapter 3 Index Structure | 19 |
| Chapter 4 Experiment | 24 |
| 4.1 Experiment Setup | 24 |
| 4.2 Experiment Result | 27 |
| Chapter 5 Conclusion and Future Work | 32 |
| Bibliography | 34 |

List of Tables

| | | |
|-----------|--|----|
| Table 1.1 | A summary of performance of various of succinct planar graph representations. Notation: n and m denote the numbers of vertexes and edges in Graph G , respectively, i is the number of isolated vertexes in G , ϵ is an arbitrary positive constant, $\tau = \min\{\lg k / \lg \lg m, \lg \lg k$, and H is the information theoretic lower bound regarding storing a separable graph. $\}$ | 5 |
| Table 4.1 | Test graphs used in our experiment | 27 |
| Table 4.2 | The performance of compact mechanisms. Time is in seconds and space is in bits per edge for encoding the edges | 28 |
| Table 4.3 | The space cost of different indexing structures. Space is in bits per edge. | 29 |
| Table 4.4 | The time performances of different indexing structures to conduct a BFS, the time cost is measured by millisecond | 30 |
| Table 4.5 | Summary of space and time performance. Space is in bits per edge and time is in millisecond | 31 |

List of Figures

| | | |
|------------|--|----|
| Figure 2.1 | An example of a partition process using edge separator | 12 |
| Figure 2.2 | An example of edge separator tree | 13 |
| Figure 3.1 | The architecture of the <i>indirect</i> indexing structure | 21 |
| Figure 4.1 | The working mechanism of the lookup table | 25 |

Abstract

This is a test document.

Acknowledgements

My Supervisor, Dr. Meng He, deserves the most praise for his invaluable guidance, friendship and suggestions. I would also like to thank all of the friends and acquaintances I have made at Dalhousie University. They have made this journey something to remember. Without everyone's support and guidance this project would have not been a success. I would like to thank Dalhousie University to provide such a great environment to do research and allowing me to finish my work. I am proudly grateful to my supportive and loving family for their whole-hearted support of my graduate studies.

Chapter 1

Introduction

Nowadays, many applications use graphs to show the relationship and represent connectivity between multiple objects. The usage of such graphs are so popular in representing various of data type including link structure of the web, 3d models in mesh, geographic maps, and surface meshes in computer graphics, etc. As the graphs inevitably grow very huge, the space issue becomes ever more important. Hence, the problem of designing space-efficient data structures to represent graphs while supporting efficient queries operations, has drawn a great deal of attention.

Plenty of previous works have been conducted to represent graphs "succinctly", which means represent the graphs in a compact representation whose space cost is close to the information-theoretic lower bound, and still supports efficient query operations. The problem of succinctly representation of a graph can be formalized as: given a graph $G = (V, E)$ of type χ , where n and m are numbers of vertexes and edges, respectively, then represent G using $\lg |\chi| + o(\lg |\chi|)$ bits of space, and support a set of query operations on the graph in constant time (assuming we can access $\Theta(\lg n)$ consecutive bits in one operation). The space bound may not be satisfied, alternatively, trade-off between space cost and query time can be further investigated.

A considerable strategy to save the space when represent graphs is taking advantage of structural properties of graphs. In practice, the most common structural property of graphs is they have small separators. A graph has separators if it can be partitioned into two subgraphs in approximately equal size by removing a small amount of vertexes. For example, Planar graphs, have $O(n^{1/2})$ separators, and even graphs that are not strictly planar because of crossings, such power networks, have small separators too. More generally, nearly all graphs that represents connections relationship in low dimensional spaces have small separators. We say a graph is separable if it is taken from a class of graphs that satisfies an n^c -separator theorem [22] for some constant $c < 1$.

In this project, we are interested in improving a compact representation mechanism for separable graphs presented by a previous work [16]. In which work the authors proposed an approach to representing separable graphs compactly. Their representation used $O(n)$ bits, meanwhile using constant time on degree or adjacency query, and neighbour listing for one vertex in constant time per neighbour (They took advantage of $O(\lg n)$ -bit parallelism computation to access $\Theta(\lg n)$ consecutive bits in one operation). Graphs with good separators got good compression by using their representation, and even graphs that are not strictly separable, their representation still works well because the separable components in those graphs can be compressed. In their paper, they provided detail description for compressing the graph by building two structures: the adjacency table and the root-find structure [16], which used vertex separators to encode the graph into a shadow adjacency table [16], as well as

support constant time on query operations. But their experiment implemented the data structure by using edge separator instead of vertex separator, which means the shadow label and root-find structure were not needed. In this project, we implement the data structure by using edge separator, and conduct experiment on it.

We firstly follow their idea to implement the data structure by using edge separator. The data structure is building by recursively partitioning a graph into two subgraphs until only one vertex left. During the partition process, an edge separator tree is built. Then make use of the edge separator tree to renumber the vertexes. We store an adjacency list for each vertex, then concatenate the adjacency lists in the order of renumbered vertex to form an adjacency table. Each pair of vertexes in the adjacency is stored by encoding the difference d between the two vertexes, and all the differences are stored contiguously as a sequence of bits in memory. We show, by using adjacency table, $O(n)$ bits are sufficient to encode the separable graph.

Next, we implement several index structures to support degree and adjacency query in constant time, and neighbour listing in constant time per Neighbour. The index structures contains all the index structures involved in their experiment, as well as two index structures implemented by us which apply succinct data structures on bit vector. In their paper, they encoded the degree of each vertex at the start of each adjacency list. However, the space of degree can be saved in some cases, which causes a trade-off between the time to support degree queries and the space to encode the graph.

We compare the performance of various indexing structure by conduct a breadth-first-search on the graph represented by different indexing structure. We also compare our space and time cost to that of an array-based adjacency list representation.

1.1 Related Works

There has been considerable works dedicate to compress graphs such as planar graph, k -Page graph and separable graph. The first attempt of succinct representation of graphs was conducted by Jacobson [11], who firstly showed how to represent planar graph by using $O(n)$ bits while supporting adjacency queries in $O(\log n)$ time. His approach decomposes a planar graph into at most four one-page graph which is represented as a sequence of balanced parentheses. His representation extended naturally to a k -page graph problem, where $k \geq 1$.

Munro and Raman [15] improved the time for adjacency queries to $O(1)$ time while using $8n + 2m$ bits where n and m are numbers of vertexes and edges respectively. Their work gives great reduction in the higher order term for number of bits regarding Jacobson's parenthesis representation, and further expands the types of operations on parenthesis representation. Geary *et al.* [13] gave a conceptually simpler representation compared to Jacobson's data structure, they used $2n + o(n)$ bits while supporting natural parenthesis operations in $O(1)$ time. The space constants on the high order term has been improved by Chuang *et al.* [21], and further by Chiang *et al.* [24]. The former proposed another encoding mechanism based on *canonical orderings* of a planar graph, then used a *multiple parentheses sequence* to represent the graph, while

the latter generalized the notion of canonical orderings to orderly spanning trees and improved constant factors in terms of numbers of vertexes and edges.

A succinct representation for a k -page graph when k is large to support various navigational operations more efficiently, was proposed by Barbay *et al.* [3], whose representation used $n + 2m \lg n + m \cdot o(\lg k) + o(m)$ bits to support adjacency queries in $O(\lg k \lg \lg k)$ time, degree queries in constant time and neighbourhood in $O(d(x) \lg \lg k)$, where $d(x)$ is the degree of vertex x , or used $n + (2 + \epsilon)m \lg k + m \cdot o(\lg k) + O(m)$ bits to support adjacency, degree and neighbourhood queries in $O(\lg k)$, $O(1)$ and $O(d(x))$ time, respectively, where ϵ is any constant satisfies $0 < \epsilon < 1$.

Table 1.1: A summary of performance of various of succinct planar graph representations. Notation: n and m denote the numbers of vertexes and edges in Graph G , respectively, i is the number of isolated vertexes in G , ϵ is an arbitrary positive constant, $\tau = \min\{\lg k / \lg \lg m, \lg \lg k\}$, and H is the information theoretic lower bound regarding storing a separable graph. }

| Type | Ref. | Space | Adjacency | Neighborhood | Degree |
|---------------|------|--|----------------------|------------------------|------------|
| Planar | [11] | $O(n)$ | $O(\lg n)$ | $O(\deg(x) \lg n)$ | $O(\lg n)$ |
| | [21] | $\frac{3}{5}m + (5 + \epsilon)n + o(n)$ | $O(1)$ | $O(\deg(x))$ | $O(1)$ |
| | [24] | $2m + 2n + o(m + n)$ | $O(1)$ | $O(\deg(x))$ | $O(1)$ |
| | [15] | $2m + 8n + o(n)$ | $O(1)$ | $O(\deg(x))$ | $O(1)$ |
| k-page | [11] | $O(kn)$ | $O(\lg n + k)$ | $O(\deg(x) \lg n)$ | $O(\lg n)$ |
| | [5] | $(2(m + i) + o(m + i)) \lg k$ | $O(\tau \lg k)$ | $O(\deg(x) \tau)$ | $O(1)$ |
| | [3] | $2m \lg k + n + o(m \lg k)$ | $O(\lg k \lg \lg k)$ | $O(\deg(x) \lg \lg k)$ | $O(1)$ |
| | [3] | $(2 + \epsilon)m \lg k + n + o(m \lg k)$ | $O(\lg k)$ | $O(\deg(x))$ | $O(1)$ |
| | [15] | $2m + 2kn + o(kn)$ | $O(k)$ | $O(\deg(x) + k)$ | $O(1)$ |
| Triangulation | [21] | $2m + n + o(n)$ | $O(1)$ | $O(\deg(x))$ | $O(1)$ |
| Separable | [16] | $O(n)$ | $O(1)$ | $O(\deg(x))$ | $O(1)$ |
| | [12] | $o(n) + H$ | $O(1)$ | $O(\deg(x))$ | $O(1)$ |

Blandford *et al.* [16] made use of the unlabelled graphs with small separators to compactly represent the graph, whose representation used $O(n)$ bits on graph satisfies

a n^c -separator theorem, while supported degree and adjacency queries in $O(1)$ time, and neighbourhood queries in constant time per neighbour. Subsequently, Blleloch *et al.* [12] proposed a succinct representation on unlabelled separable graphs, which used $lg|\chi| + o(n)$ bits to represent a graph in type χ , while holding the same time bounds as the performance of Blandford's representation. A summary of performance of various of succinct planar graph representations is shown in Table .

1.2 Preliminaries

Graph Separator. A family of graphs G is define separable if : it is closed under taking subgraphs, and satisfies the $f(\cdot)$ -separator theorem [22] if there is a constant $\alpha < 1$ and $\beta > 0$ such that each member graph in G with n vertexes has a separator set S of size $|S| < \beta f(n)$, that partition the graph into two parts A and B , with at most αn vertexes in each. A graph is separable if it belongs to a separable family of graphs.

In this project we focus on the graphs that satisfy the n^c -separator theorem for some constant $c < 1$. One class of graphs that reach this specifications is the planar graphs, which satisfies the theorem that $c = \frac{1}{2}$. Another example can be well-shaped meshes in \mathbb{R}^d , with separators of size $O(n^{1-1/d})$ [18].

Lipton *etal.* [17] prove that all classes of graph which satisfy $n/(\log n)^{1+\epsilon}$ -separator theorem have bounded density. The bounded density means that every n -vertex member in a class if graphs has bounded $O(n)$ edges. So we can assert that the

separable graphs have bounded density.

By making use of the definitions above, we define a class of graphs G satisfies a $f(n)$ -edge separator theorem if there are constant variables $\alpha < 1$ and $\beta > 0$ such that each member graph in G with n vertexes, has an edge separator with at most $\beta f(n)$ edges whose removal partitions the graph into two subgraphs with at most αn vertexes in each. The edge separator is not as common as the vertex separator, because a graph with an edge separator of size s also has a vertex separator of size s at most, but no similar bounds holds when it is conversed [16].

In this project, we will only consider the connected graph, which means all vertexes in the graph have nonzero degree, and we assume each step of partitioning always returns a edge separator set of size $O(n^c)$.

Queries. Our data structure supports three kinds of queries on separable graphs: adjacency queries, degree queries and neighbour listing. The adjacency query checks whether there is an edge between two vertices, the neighbour listing returns the neighbours of a given vertex, and the degree query reports the number of edges connected to a given vertex.

Bit Vector. The compact representation stores a separable graph as difference-encoded bit sequence in memory. To support the three kinds queries (adjacency queries, degree queries and neighbour listing) in constant time, knowing both the

number of encoded adjacency lists in the sequence up to an index position, and the start position of a particular adjacency is necessary. More formally, the compact representation requires data structures that support rank and select operations in constant time: Given a set $B[0...n)$ to represent a subset S of a universe $U = [0...n) = \{0, 1, \dots, n-1\}$, where $B[i] = 1$ *iff* $i \in S$,

- $rank_q(x) = \{k \in [0...x] : B[k] = q\}$
- $select_q(x) = \min\{k \in [0...n) : rank_q(k) = x\}$

In this project, we use two kinds of succinct data structure on bit vector : RRR [20] and SD vector [19], to build the index structure which supports queries on separable graph. Both structures satisfy information-theoretic lower bound, and still supports efficient rank and select query operations.

Adjacency Tables. An adjacency list representation for a graph associates each vertex in the graph with the collection of its neighbouring vertices. For use as a data structure, the main alternative to the adjacency list is the adjacency matrix. For a sparse graph (most pairs of vertices are not connected by edges), in our case, a graph has *bounded mdensity*, an adjacency list representation is significantly more space-efficient than an adjacency matrix representation: the space cost of the adjacency list representation depends on the number of edges and vertices in the graph, while the adjacency matrix representation is stored as the square of the number of vertices. Generally, an adjacency list representation for an undirected graph requires $8|E|$ bytes

$(2 \cdot 32 \cdot |E|)$ bits) of space, where $|E|$ is the number of edges in the graph. However, it is not space-efficient enough yet.

In our project, we assign each vertex in the graph an integer label. For each vertex, we store an adjacency list which contains its all neighbouring vertexes. If a vertex with label v has neighbours $v_1, v_2, v_3, \dots, v_n$ in ascending sorted order, then we encode the difference between each adjacency pair in the adjacency list as $|v_1 - v|, |v_2 - v_1|, |v_3 - v_2|, \dots, |v_n - v_{n-1}|$ contiguously, which forms a sequence of bits stored in memory. We use the Elias gamma code [9] to encode the differences, which uses $2\lfloor \log n \rfloor + 1$ bits to encode a difference. The difference $v_1 - v$ may be negative in some cases, we here store a one-bit flag that value. To implement one of our index structures, the degree of each vertex needs to be encoded at start position of the corresponding adjacency, but in other circumstances the space of degree can be saved. By concatenating all the adjacency lists in the order of the vertex labels, we form an adjacency table. To access the adjacency list of a vertex, knowing the starting position of the adjacency is necessary. The easiest and fastest approach is to store an array of offset pointers of size $O(\log(n))$ for each. But if doing so, we would use $nO(\log(n))$ space, which exceeds our space bound. Alternatively, we implement other two index structures proposed by Blandfor [16], as well as two index structures based on succinct data structure on bit vector, which use less space cost down to $O(n)$ bits.

Lemma 1.2.1. *An adjacency table can support degree queries in $O(1)$ time, and*

neighbour listing in $O(d)$ time, where d is the degree of the vertex.

Proof. To access an adjacency list of a particular vertex, we use our index structure to locate its starting position in constant time. Each adjacency list consists of a encoded degree number and a sequence of encoded differences, which use $O(\log(n))$ bits and $dO(\log(n))$ bits respectively. By taking advantages of the $O(\log n)$ -bits parallelism computation, we can decode any $O(\log(n))$ -bits value in constant time, like using lookup table. Hence we can decode the degree at the starting of each adjacency list in constant time. For neighbour listing, we decode each difference pair in constant time. So it takes constant time to make degree queries, and $O(d)$ time to make neighbourhood queries, where d is the degree of the vertex. \square

1.3 Outline

This report is organized as follows. Chapter 1 shows the motivation of this project, and provides some preliminaries used in this project. Chapter 2 and 3 explain the construction and usage of the edge separator tree and several index structures respectively. Chapter 4 describes the experiment setup and examines the performance of our data structure. Chapter 5 shows the conclusion as well as the future work of our project.

Chapter 2

Edge Separator Tree

To compactly represent the tree, the first thing we need to do is building an edge separator tree from the original tree. The tree-building process is based on recursively partition the tree into two parts using edge separator. One example that illustrates the function of edge separator is shown in Figure 2.1: assuming we have an 8-vertex graph G_1 with 9 edges, by removing the two edges (v_1, v_9) and (v_1, v_5) , we can partition G_1 into two parts with 4 vertexes in each. Hence the edge separator in this partitioning step is $\{(v_1, v_9), (v_1, v_5)\}$.

We can see that each time we make a partitioning on a graph, it splits its vertex set and parts of edges set, and removes some edges. if we perform this partitioning recursively until only one vertex left, then we have an edge separator tree from the graph, one edge separator tree for the G_1 in Figure 2.1 is shown in Figure 2.2. Each vertex in a graph will appear once in a leaf of the edge separator tree build from that graph. Each internal node is the set of edges used as edge separators to make partitioning in each step.

Note that we arbitrarily decide which side of partition will be the left child or the right child in this stage. The tree is also not perfect balanced in some cases (e.g.

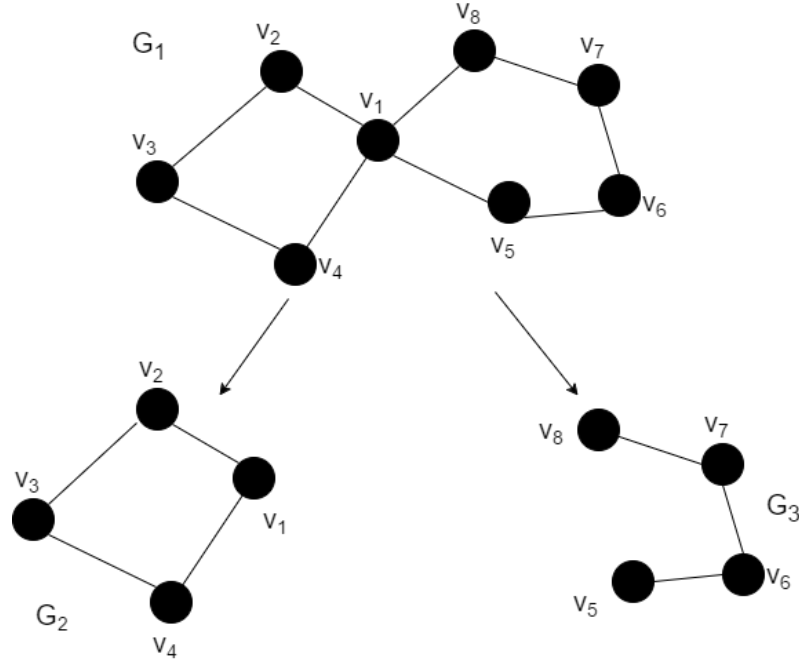


Figure 2.1: An example of a partition process using edge separator

one child is a leaf and the other is an internal node), the side of the children is also decided arbitrarily. We will take advantages of this "freedom" to further reduce the space to encode the graph in later time. The vertex labels may be disordered by making in-order traversal on the edge separator tree (e.g. in Figure 2.2 the order of the vertexes after in-order traversal is $:v_1, v_4, v_2, v_3, v_5, v_6, v_7, v_8$), we will renumber these vertex by making a in-order traversal along the edge separator tree.

METIS Partition Library. To recursively partition the graph, we use METIS's API in our project. METIS is an open source library developed by Karypis lib in University of Minnesota. It supports set of serial programs for partitioning graphs and partitioning finite element meshes. The algorithms implemented in METIS are based on the multilevel recursive-bisection, multilevel k-way, and multi-constraint

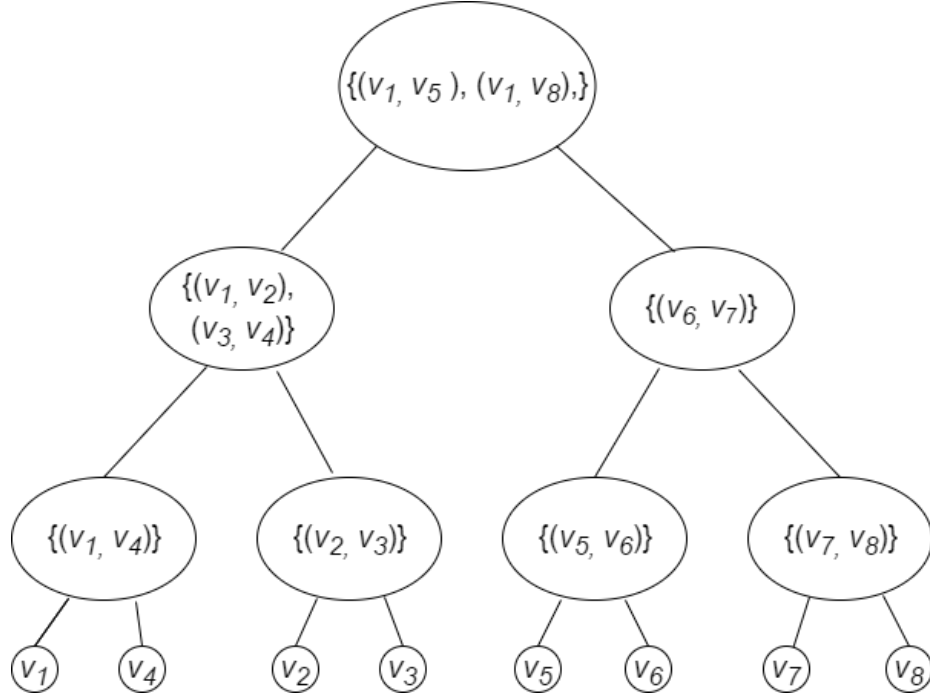


Figure 2.2: An example of edge separator tree

partitioning schemes [2]. The METIS partitioning API returns one integer list with vertex index as key, and a flag 0 or 1 to indicate which side of partitioning the vertex belongs to. The returned value of METIS API does not provide information about edges, and the API takes two integer arrays (adjacency list of current graph and the index of vertex in adjacency list) as input parameters, hence before each calling of the METIS API, we need to generate an adjacency list and an index list for the current graph to be partitioned, and in order to collect the edge separators, we need both the adjacency list and the returned values after calling the API.

The Algorithm to build the separator tree is given in Algorithm 1. Here we use (V, E) to represent a graph for simplicity, but in fact the METIS API takes formalized adjacency list and the index list as input parameters. So each time to call the METIS

API, the current vertex set need to be renumbered temporally in contiguous ascending numbers, and the adjacency list needs to be initialized according to the temporally order. The METIS API also does not support to partition the graph with the vertex amount under 3, hence the very basic partitioning cases are defined by ourself.

Algorithm 1: Building the edge separator tree

```

1 buildTree ( $V, E$ );
2 if  $|V| = 1$  then
3   | return  $V$ ;
4 else
5   |  $(V_a, V_b) \leftarrow \text{METIS\_API} (V, E)$  ;
6   |  $E_{sep} \leftarrow \text{findEdgeSeparator} (V, E, V_a, V_b)$  ;
7   |  $E' = E - E_{sep}$  ;
8   |  $E_a \leftarrow \{(u, v) \in E' \mid u \in V_a \vee v \in V_b\}$  ;
9   |  $E_b = E' - E_a$  ;
10  |  $T_a \leftarrow \text{buildTree} (V_a, E_a)$  ;
11  |  $T_b \leftarrow \text{buildTree} (V_b, E_b)$  ;
12  | return separatorTree ( $T_a, E_{sep}, T_b$ ) ;
13 end
```

Building Adjacency Table. After getting the edge separator tree, we need to renumber the vertexes based on an in-order traversal along the tree. Each vertex appears in one leaf of the tree will receive a new label which increases by one when visiting a new leaf node. A mapping table is also generated after completing the traversal, which will be used to build an adjacency list for each renumbered vertex by using the new vertex labels. Then we can build an adjacency list for each vertex, the list contains all the neighbouring vertexes in sorted order. We use gamma code to encode the difference of each adjacency pair in the list (e.g. $v_1 - v$, $v_2 - v_1, \dots$, $v_{n-1} - v_{n-2}$), each difference d costs $2\lceil \log n \rceil + 1$ bits to be encoded. Note that the

value $v - v_1$ might be negative, so we store one flag bit for it. We also encode the number of the entries in an adjacency list at the starting position of the list, which will be used to support degree queries in constant time.

Lemma 2.0.1. *Any n -vertex member in a class of graphs which satisfies a n^c – edge separator theorem, can be encoded in $O(n)$ bits using an adjacency table.*

Proof. For each edge (u, v) in an edge separator in a graph with s vertexes, the difference between vertex u and vertex v is $O(s)$. We use gamma code to encode the difference which would use $O(\log(s))$ bits, so that edge will contribute $O(\log(s))$ bits to the adjacency lists of both vertex u and vertex v . Hence, we charge $O(\log(s))$ to every edge in a separator of a graph with s vertices. Recall that by using the n^c – edge separator theorem, each member with n vertexes in the class of the graph has a βn^c edge separator whose removal partitions the member into two parts with at most βn vertexes in each. Let define $S(n)$ to be an upper bound of required bits to encode a graph with n vertexes, and let $\alpha < a < 1 - \alpha$, then $S(n)$ satisfies the recurrence:

$$S(n) \leq S(an) + S(n - an) + O(n^c \log n)$$

This recurrence can be solved to $S(n) = O(n)$, and the bits to encode the degree of all the vertexes is bounded by $O(n)$, because there are $O(n)$ edges in total. So the total space is $O(n)$ bits. □

Child-flipping. When encoding the graph, We contiguously encode the difference

between each adjacency pair in the list by using gamma code. The gamma code use $2\lfloor \log n \rfloor + 1$ bits to encode a difference d . If we can make each difference "smaller", in other words, if we can make the labels of vertexes in an adjacency list closer to each other, then we can reduce the space to encode the adjacency list. We arbitrarily decide the side of partitioning during the construction of the separator tree (in our project, we decide the side of partitioning according to the output value of METIS API: '0' to the left child and '1' to right child), so there is a degree of freedom in the way to build the tree. To take advantage of this freedom, we can apply a heuristic optimization method called "child-flipping" [16] to our edge separator tree.

The child-flipping algorithm works when making traversal along a separator tree, it tracks the nodes containing vertex labels which appear before and after the vertex labels in the current node. Let define N_L to be the left children of current node's left ancestors and N_R to be the right children of the current node's right ancestors, N_1 and N_2 represents the current node's left child and right child, respectively, and E_{AB} indicates the number of edges between the vertexes in node A and node B. For each node, the vertexes in N_1 of the current node will have labels less than the vertexes' under the node, the vertexes in N_2 will have larger labels than all the vertexes under the node, and there is no intersection between the vertexes in N_L and N_R of the node. The child-flipping exams $E_{N_L N_1}$, $E_{N_2 N_R}$, $E_{N_L N_2}$ and $E_{N_1 N_R}$ of each node to ensure that $E_{N_L N_1} + E_{N_L N_2} \leq E_{N_L N_1} + E_{N_L N_2}$. If not, the child-flipping algorithm swaps the two children of current node. The meaning of making this swap is : after renumbering the vertexes based on an in-order traversal, if $E_{N_L N_1} + E_{N_L N_2} < E_{N_L N_1} + E_{N_L N_2}$, then

the space to encode the edges $N_L N_2$ and $N_1 N_R$ is greater than that to encode $N_L N_2$ and $N_1 N_R$, because the difference $v_{N_2} - v_{N_L}$ is obviously greater than $v_{N_1} - v_{N_L}$, similar fact works on the difference of $v_{N_R} - v_{N_1}$ and $v_{N_R} - v_{N_2}$.

This heuristic algorithm can be applied to any edge separator tree as a postprocessing step. However, this process is quite time-consuming, the time cost of applying both METIS and child-flipping in Blandfor *etal.* [16]'s experiment shows a demonstration on this. To get the four numbers of edges ($E_{N_L N_1}$, $E_{N_2 N_R}$, $E_{N_L N_2}$ and $E_{N_1 N_R}$), we need to have the four vertex sets: N_L , N_R , N_1 and N_2 for each tree node. The N_L and N_R of each node in separator tree can be pre-collected, but the N_1 and N_2 need to be calculate locally when reach the node. In our project, we make some efforts to reduce the time to conduct this child-flipping processing: Let assume we have a node N_a in a separator tree with left child as N_b and right child as N_c , if we knew the N_1 and N_2 of N_a (denote as N_{1a} and N_{2a}), then

- for node b : $N_{1b} = N_{1a}$, $N_{2b} = N_{2a} \cup N_{Ra}$
- for node c : $N_{1c} = N_{1a} \cup N_{La}$, $N_{2c} = N_{2a}$

Having this general case, then we can generate the N_1 and N_2 for each node during the child-flipping traversal via DFS. The N_L and N_R of each node in separator tree can be pre-collected via a post-order traversal along the edge separator tree, the vertex set of a node is formed by merging the two vertex sets of its children in sorted order. When needing to expand the current N_1 or N_2 , we include the new coming N_L or N_R as distinct sorted set instead of merging N_L into N_1 or N_R into N_2 . Note that

in term of a node, the size of N_L or N_R is always smaller than the size of N_1 or N_2 , respectively. We check an existing of an edge by performing binary on larger sorted set, all of above approaches improve the speed of child-flipping processing significantly

.

Chapter 3

Index Structure

We already introduce how to encode a graph with good edge separator into difference adjacency table bounded by $O(n)$ bits, but to access to the particular adjacency list for a particular vertex, we need to know the exact starting position in the bit sequence. Hence our representation requires a data structure that support efficient select operations on the bit sequence. Note that, the space cost of the data structure should not exceed our space bound which is $O(n)$. Blandfor *etal.* has proposed three kinds of index structure: *direct*, *semi-direct*, and *indirect* [16] to support efficient select queries. In addition to these three structures, we also implement other two index structures based on RRR and SD Vector respectively, we will henceforth call these two RRR-based and SDV-based index structure.

Direct indexing structure. The easiest way to access to the starting position of any adjacency list is to store an offset pointer for each vertex. Because our representation uses $O(n)$ bits to store a separable graph, one offset pointer will use $\Theta(\log(n))$ bits. For all vertexes, it will use $n\Theta(\log(n))$ in total. If we want to locate the starting position of any vertex, only one memory access is needed, so this approach is fastest among the five index structures. In our implementation, we allocate one word (32

bits) space for each vertex, hence the total space for this direct index structure is $4n$ bytes.

Semi-direct indexing structure. This *semi – direct* index structure is very similar to the *direct* structure, but instead of allocating one word space for each vertex, we use two words to store four offset values. Let say we have four vertexes $v_n, v_{n+1}, v_{n+2}, v_{n+3}$, The first word space will store the pointer value of v_n , then use the second word to store three offset value $v_{n+1}-v_n, v_{n+2}-v_n, v_{n+3}-v_n$. If the space of the three offset values exceeds one-word bound, then store the offsets somewhere else, then the second word store a pointer pointed to them. Compared to the *direct* structure, if the vertex satisfies $v_{id} \bmod 4 = 0$, then it works exactly as same as *direct* structure, if not, it firstly gets the real offset by adding the corresponding offset to the first pointer value, then access the memory directly. If every four vertexes fit into two-word space, then this *semi – direct* structure would use half of the space used by *direct* structure.

Indirect indexing structure. Note that the above two index structure use $n\theta(n)$, which exceeds $O(n)$ space bound. To limit the space cost within the bound, Blandfor *etal.* implement the third structure : *indirect*, which use $O(n)$ space meanwhile supports constant time on locating vertex. The component of the structure is shown in Figure 3.1: we firstly divide all vertexes into blocks with $\log(n)$ vertexes in each, then we further divide each block into subblocks. Each subblock contains a minimal number of vertexes whose indexing offset length reaches at least $k \log(n)$

bits for some constant number k . We store a bit vector with $\log(n)$ bits length for each block, if v_i is the first vertex in its subblock, then the $(i \bmod \log(n))_{th}$ bit in the vector of its block is set to 1, otherwise 0. We only store the offset pointers of vertexes if their corresponding bits in bit vector are 1. This will require $O(n)$ bits in total.

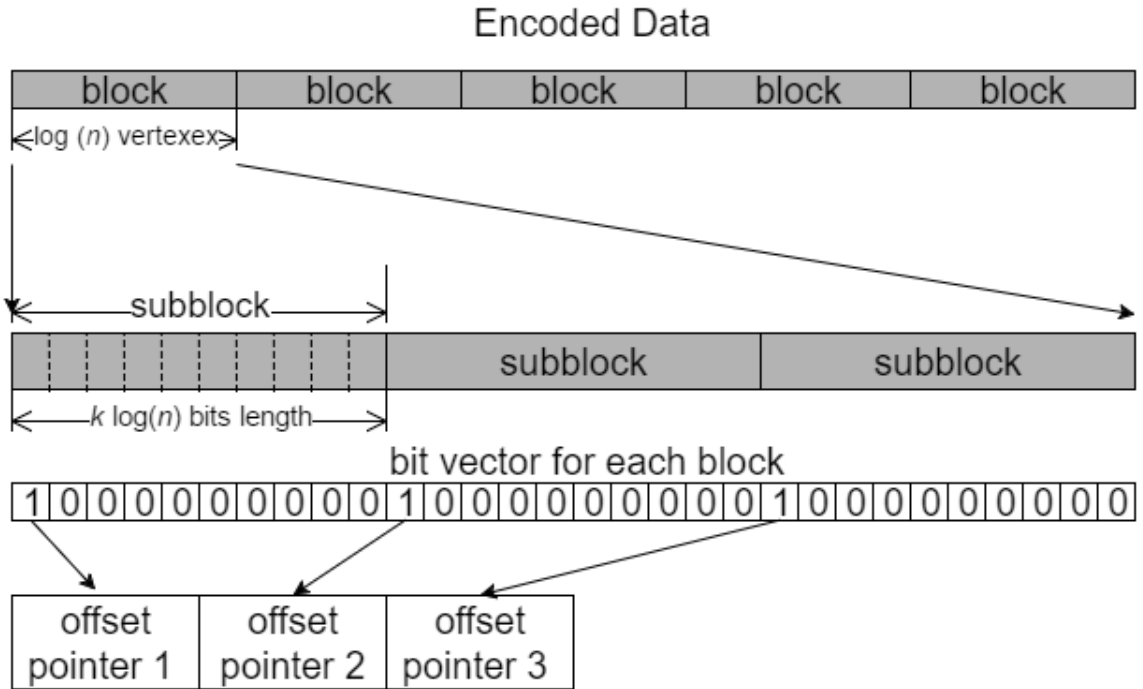


Figure 3.1: The architecture of the *indirect* indexing structure

To locate the starting position of a particular vertex, we firstly find which block the vertex is in by dividing the vertex label with $\log(n)$ (number of vertexes in each block). Next we check the bit vector of that block to find which subblock contains the vertex. Since we store the offset pointer for the vertexes which are the first one in their subblock, we can know the starting position of the target subblock, then we decode the subblock. Although we will get all the adjacency lists within the subblock

if do so, the degree information of each vertex is also encoded, hence we can have the adjacency list of the target vertex easily. The bit vector of each block has $\log(n)$ bits, each subblock is $k \log(n)$ bits long, determining and decoding the subblock can be completed by using lookup table, which uses constant time to check $\Theta(\log(n))$ bits.

RRR-based indexing structure. RRR is a succinct data structure which was first proposed by Raman *etal.* [20]. This data structure works on bit sequence $S[0, \dots, n)$ in such a way that provides rank and select operations in $O(1)$ time at the price of $nH_0(S) + o(n)$ space cost, where $H_0(S)$ is the empirical zero-order entropy of the sequence. The RRR data structure divides the original bit sequence into several superblocks, then divide the superblocks further. For each block it stores a number to indicate the number of ones or zeros in the block. The number would be used as a lookup key in a table. Besides that, it also stores an offset which indicates which of the blocks the number is in. Offset is also used as an index in the table. By grouping blocks into superblock, RRR can avoid iterating over each block to conduct a rank or select query.

RRR efficiently compresses a bit set if the set is very sparse populated, we take advantage of this feature of RRR and implement a indexing structure based on it. The original graph is encoded into a bit sequence by using gamma code, and we know the length of this sequence. Hence we can maintain another bit set with same length as the encoded bit sequence, in which we set a bit to '1' if it is the starting point of an adjacency list, otherwise we set the bit to 0. By doing so, we get a bit set

which is very '1' sparse, then we can build RRR bit vector on this bit set. To locate a particular vertex, we only need to conduct a *select* query on the RRR bit vector, and we will know its starting position in the encoded bit sequence in constant time. Since the encoded bit sequence is $O(n)$ bits long, and for a n -bit sequence, RRR needs $nH_0(S) + o(n)$ bits to support constant *rank* and *select* queries. We do not need to store additional information besides RRR structure, so the space cost will not exceed our bound.

SD Vector-based indexing structure. SD Vector is a bit vector that can compress very sparse populated bit vectors by representing the positions of 1 by the Elias [8] - Fano [10] representation for non-decreasing sequences. Given a bit set $B[0, \dots, n-1]$ with m ones ($m \ll n$), the SD Vector uses $m \log n/m + 2m + o(m)$ [7] bits to store the bit set, meanwhile support constant time *rank* and *select* queries. The idea of SD Vector-based indexing structure is very similar to RRR-based indexing structure. We firstly maintain a bit sequence which has the same length as the encoded graph bit sequence, next we mark the bits which indicate the starting position of an adjacency list, then we build a SD Vector on this bit sequence. In order to locate the starting position of a particular vertex, we use the SD Vector to perform a *select* operation which would return the expected position.

Chapter 4

Experiment

After implementing our data structure by encoding the graph into an adjacency table as well as building an auxiliary indexing structure for supporting queries on graph, we'd like to exam the performance of our approach. In this chapter we will describe the experimental setup, then show and discuss our experiment result.

4.1 Experiment Setup

RRR and SD Vector. As we introduced in Chapter 3, we take advantage of the features of succinct representation of bit vector to build two indexing structures: RRR-based and SD Vector-based. We implemented both RRR and SD Vector structures by using the succinct data structure API provided by Simon Gog [14], which presents a various of succinct representations including bit vector, Integer Vectors, Wavelet Trees, Compressed Suffix Arrays (CSA), etc. The API for RRR and SD Vector works when a bit set is provided, then it builds *select_support_type* or *rank_support_type* data structure on it. A *select* or *rank* operation then can be called by using member functions within the *select_support_type* or *rank_support_type* objects.

Decoding. Taking advantage of the feature of Random-Access-Machine that

costs constant-time to perform operations on $O(\lg n)$ bits, we can decode any $O(\lg n)$ -bit words in constant time. One practical approach here is using a lookup table. In our project, the graph is Elias gamma encoded, to ensure a decoding operation on $O(\lg n)$ -bits requires constant time, we implement a table to decode multiple gamma codes at once. For a fixed number of k_g , we create a table of size 2^{k_g} , and use the numbers range from 0 2^{k_g} as indexes. Each entry in the table indicates a decoding result which is represented by a 32-bit word.

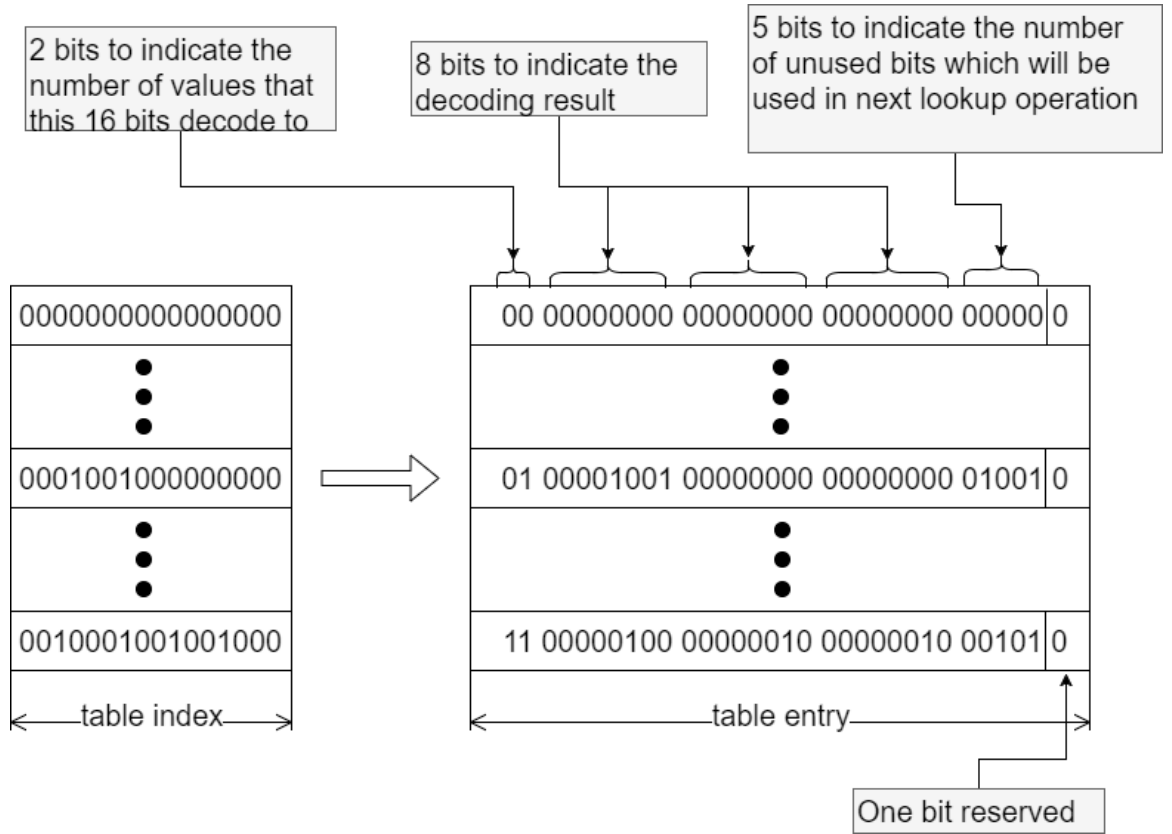


Figure 4.1: The working mechanism of the lookup table

The working mechanism of this lookup table is shown in Figure 4.1 (note that we set the k_g to 16, and limit the number of words decoded at once to 3 in our

implementation): the first two bits indicate the number of values that this 16 bits decode to. If a bit block with size of 16 can not be directly decoded (e.g. the first case in the Figure 4.1, which is all zero), then the firstly two bits in the corresponding entry is set to 0, which means none of number can be decoded in this block. The next 24 bits ($3_{rd} - 26_{th}$ bits) indicate the three potential decoding result (if the first two bits show that only one number can be decoded, then we only look into the $3_{rd} - 11_{rd}$ bits). The next 5 bits ($27_{th} - 31_{th}$ bits) show the number of unused bits which will be used in next lookup operation, and the last bit is unused. This makes each table entry fit into 32 bits. If the first two bits indicate that none value can be directed decoded by checking this table entry, then the sequence will be decoded explicitly (counting zeros then get the value).

Test Graphs. The test graph used in our experiment come from several source: two 3D Mesh graphs from an online Graph Partitioning Archive maintained by C. Walshaw [6], one graph from Scotch Graph Archive maintained by Dave Beckett [4], one graph from RIACS Graph Archive [2], two Gpu Layout Plugin graphs from Cytoscape online resource [23], and two BCS Structural Engineering Matrix graphs from Matrix Market [1]. All the test graphs are undirected graphs, a summary of general attributes of the graphs is shown in Table 4.1, giving the information about the numbers of vertexes and edges, max degree and graph source.

Table 4.1: Test graphs used in our experiment

| Graph | Vertexes | Edges | Max Degree | Source |
|----------|----------|---------|------------|-------------------|
| feocean | 143437 | 409593 | 6 | 3D mesh [6] |
| m14b | 214765 | 1679018 | 40 | 3D mesh [6] |
| febody | 45087 | 163734 | 28 | Scotch Graphs [4] |
| wave | 156317 | 1059331 | 44 | RIACS Grids [2] |
| 598a | 110971 | 741934 | 26 | Cytoscape [23] |
| 144 | 144649 | 1074393 | 26 | Cytoscape [23] |
| bcsstk30 | 28924 | 1007284 | 218 | BCS [1] |
| bcsstk31 | 35588 | 572914 | 188 | BCS [1] |

4.2 Experiment Result

. The experiments were conducted on a Linux server containing 32 Intel Xeon E5-2650 @ 2.00GHz cores and 256 Gigabytes main memory. We present each edge in undirected graphs in both directions, which means the space cost measured by bits per edge is reported for each direction. Our experiment focus on three parameters of concern: the time to construct our data structure, the time to perform queries on graph and the space needed to encode the graph into adjacency table. The space cost of our data structure consists of the space for adjacency table, and the space for the indexing structure. There are a trade-off between the time to complete the construction and the space for encoding the adjacency table caused by whether conduct the child-flipping process, as well as a trade-off between the space for the indexing data structure and the time to perform queries on graphs resulted by applying different indexing structures. Our experiments present demonstration on these trade-offs.

Table 4.2 shows the time cost of constructing the compact data structure and the space needed to encoding the adjacency table. We can see that the child-flipping

Table 4.2: The performance of compact mechanisms. Time is in seconds and space is in bits per edge for encoding the edges

| Graph | METIS | | METIS-CF | | Degree |
|----------|-------|-------|----------|-------|--------|
| | Time | Space | Time | Space | |
| feocean | 55.69 | 8.68 | 59.51 | 8.27 | 0.86 |
| m14b | 129.3 | 5.21 | 155.37 | 5.00 | 0.51 |
| febody | 6.33 | 4.64 | 7.2 | 4.25 | 0.82 |
| wave | 77.4 | 5.76 | 86.52 | 5.57 | 0.55 |
| 598a | 39.06 | 5.45 | 52.33 | 5.21 | 0.54 |
| 144 | 67.09 | 5.32 | 84.33 | 5.13 | 0.52 |
| bcsstk30 | 7.89 | 2.29 | 12.3 | 2.24 | 0.17 |
| bcsstk31 | 6.61 | 2.98 | 9.02 | 2.89 | 0.30 |

causes additional time cost as a heuristic postprocessing phase, but do reduce the space cost of the adjacency table (save at most 8% space on the febody graph, and at least 2.5% space on the bcsstk30 graph). The degree space for each graph is listed separately, because the child-flipping process will not change the degree of each vertex, and the number of bits to encode a number is fixed (use $2\lfloor \log d \rfloor + 1$ bits to encode a degree d by using Elias gamma coding).

The indexing structures are built based on the ordering result after the child-flipping process. Table 4.3 illustrates the space cost of each indexing structure, while Table 4.4 illustrates the time cost to perform a breadth-first-search(BFS) on the graph by applying those indexing structures on the encoded adjacency table. The reason why we use BFS to measure the performance is because it requires visiting all the edges, which makes it a reasonable measure for our data structure. The BFS algorithm starts at some arbitrary vertex as the root (we use the 1_{st} vertex after renumbering in our project) and explores the neighbour nodes first, before moving to the next level neighbours, with a cost of time $\Theta(|V| + |E|)$.

We also compare the performance of our representation to that of an array-based adjacency list graph representation. The array-based adjacency list representation store all the neighbours of each vertex contiguously in a large array with the neighbour lists place one after the other. It also require an array to store the index of each vertex. Therefore, it use one 32-bit word to present an edge and one 32-bit word as the index of one vertex, taking $8|E|$ bytes ($2 \times 32 \times |E|$ bits) space for list and $4|V|$ bytes for indexes in total. For the BFS algorithm on any types of indexing structures, we use one bit flag to indicate whether a vertex has been visited.

Table 4.3: The space cost of different indexing structures. Space is in bits per edge.

| Graph | Direct | Semi | Indirect | RRR | SD Vector |
|----------|--------|------|----------|------|-----------|
| feocean | 5.63 | 2.81 | 0.84 | 2.03 | 1.59 |
| m14b | 2.05 | 1.02 | 0.42 | 1.01 | 0.64 |
| febody | 4.46 | 2.23 | 0.54 | 1.33 | 1.31 |
| wave | 2.36 | 1.18 | 0.46 | 1.14 | 0.72 |
| 598a | 2.39 | 1.19 | 0.47 | 1.10 | 0.77 |
| 144 | 2.15 | 1.07 | 0.43 | 1.05 | 0.67 |
| bcsstk30 | 0.46 | 0.23 | 0.12 | 0.35 | 0.16 |
| bcsstk31 | 0.93 | 0.46 | 0.21 | 0.54 | 0.35 |

The experimental result shows that the *direct* and *semidirect* indexing structures give the fastest query time, but still slower than the array-based adjacency list representation. It's not surprising because the array-based representation has good spacial locality, which means the edges of a vertex are adjacent in memory and be loaded into cache as an array. In addition to that, the array-based representation does not need to decode the list, hence the array-based representation is faster than others. The *semi – direct* indexing structure uses half of the space when compared to *direct* indexing structure while requiring little extra time on query. *Indirect* indexing

structure shows good space usage but it is not very efficient in query, because it needs to decode a "subblock" rather than an adjacency list. The SD Vector-based structure outperforms the RRR-based structure in both space usage and time cost in all the test graphs, but since we are using the API provided by a third-party to construct RRR and SD Vector on bit set, we can not simply say the RRR structure show worse performance than SD Vector when applied on our indexing structures.

Table 4.4: The time performances of different indexing structures to conduct a BFS, the time cost is measured by millisecond

| Graph | Array | Direct | Semi | Indirect | RRR | SD Vector |
|----------|-------|--------|------|----------|-------|-----------|
| feocean | 13.6 | 31.6 | 34.9 | 97.2 | 103.3 | 60.8 |
| m14b | 36.2 | 88.5 | 93.1 | 216.8 | 251.6 | 130.7 |
| febody | 2.5 | 7.3 | 7.9 | 27.7 | 33.7 | 13.8 |
| wave | 19.9 | 60.8 | 66.1 | 144.8 | 176.2 | 95.1 |
| 598a | 14.5 | 42.1 | 45.9 | 102.8 | 115.0 | 65.6 |
| 144 | 19.2 | 60.0 | 64.7 | 111.2 | 130.4 | 87.9 |
| bcsstk30 | 4.2 | 24.9 | 26.8 | 37.0 | 40.9 | 29.3 |
| bcsstk31 | 3.7 | 20.2 | 20.7 | 35.4 | 39.2 | 27.9 |

Table 4.5 illustrates a summary of performance of different indexing structures associated with the encoded adjacency table built after child-flipping process. The space for each mechanism contains the space for the encoded adjacency list, the space for degree and the space for the indexing structure. We can see that the *direct* indexing structure gives fastest query speed while requiring the most space, *indirect* indexing structure use least space while needing more time on query, the SD Vector-based indexing structure presents a good trade-off between the space usage and query efficiency.

Table 4.5: Summary of space and time performance. Space is in bits per edge and time is in millisecond

| Graph | Array | | Metis/Di | | Metis/SemD | | Metis/Ind | | Metis/RRR | | Metis/SD_V | |
|-------|-------|------|----------|------|------------|-------|-----------|-----|-----------|-------|------------|------|
| | T | S | T | S | T | S | T | S | T | S | T | S |
| feo* | 13.6 | 37.6 | 31.6 | 14.8 | 34.9 | 11.94 | 97.2 | 9.8 | 103.3 | 11.06 | 60.8 | 10.6 |
| m14b | 36.2 | 34.1 | 88.5 | 7.6 | 93.1 | 6.53 | 216.8 | 5.9 | 251.6 | 6.52 | 130.7 | 6.1 |
| feb* | 2.5 | 36.4 | 7.3 | 9.5 | 7.9 | 7.3 | 27.7 | 5.6 | 33.7 | 6.4 | 13.8 | 6.4 |
| wave | 19.9 | 34.3 | 60.8 | 8.5 | 66.1 | 7.3 | 144.8 | 6.5 | 176.2 | 7.2 | 95.1 | 6.8 |
| 598a | 14.5 | 34.4 | 42.1 | 8.1 | 45.9 | 6.9 | 102.8 | 6.2 | 115.0 | 6.8 | 65.6 | 6.5 |
| 144 | 19.2 | 34.1 | 60.0 | 7.8 | 64.7 | 6.7 | 111.2 | 6.1 | 130.4 | 6.7 | 87.9 | 6.3 |
| b30* | 4.2 | 32.4 | 24.9 | 2.8 | 26.8 | 2.6 | 37.0 | 2.5 | 40.9 | 2.7 | 29.3 | 2.5 |
| b31* | 3.7 | 33.0 | 20.7 | 4.2 | 20.7 | 3.6 | 35.4 | 3.4 | 39.2 | 3.7 | 27.9 | 3.5 |

Chapter 5

Conclusion and Future Work

In this project, we focus on how to represent a separable graph compactly while supporting efficient query. Our theoretical basis comes from a previous work conducted by Blandfor *et al.* [16], which provides the theory of using vertex separator to compactly represent the graph. The representation takes $O(n)$ bits, meanwhile using constant time on degree or adjacency query, and neighbour listing for one vertex in constant time per neighbour.

We first make a survey on the researches related to succinct data structure on planar graph, as well as some succinct representation on bit vector like RRR. Then focus on the representation of separable graph.

Next, we use the METIS_API to partition the graph by removing some amount of edges recursively, and construct the edge separator tree for the graph. To reduce the space of difference encoding, we conduct a heuristic postprocessing phase called "child-flipping" on the edge separator tree. We use Elias gamma code to encode the adjacency list of each vertex, and concatenate all the adjacency lists to an adjacency table.

Then we implement several indexing structures to support adjacency queries, degree queries and neighbourhood queries in constant time. Besides the three indexing structure mentioned in Blandford's work, we implement other two structure based on succinct representation on bit vector.

Finally we conduct some experience to measure the performance of our representation. The experiments demonstrates the trade-off between the space cost to encode the adjacency table and the time to construct the data structure, as well as the trade-off between the space cost and time of queries by different indexing structures. Compared to the array-based adjacency list representation, our approach reduces space usage by almost an order of magnitude, while supporting depth-first traversal in about the same running time.

Bibliography

- [1] *BCS Structural Engineering Matrix*. Available at <http://math.nist.gov/MatrixMarket/data/Harwell-Boeing/bcsstruc5/>.
- [2] *RIACS GRIDS*. Available at <http://riacs.edu/pub/grids/>.
- [3] He M Barbay J, Castelli Aleardi L and Munro JI. Succinct representation of labeled graphs. *Algorithmica* 62, pages 224–257, 2012.
- [4] Dave Beckett. *Internet Parallel Computing Archive*. Available at <http://wotug.org/parallel/libraries/communication/scotch/Graphs/>.
- [5] Gavaille C and Hanusse N. On compact encoding of pagenumber. *Discret Math Theor Comput Sci* 10(3), 2008.
- [6] C.Walshaw. *Graph partitioning archive*. Available at <http://www.gre.ac.uk/~c.walshaw/partition/>.
- [7] K. Sadakane D. Okanohara. Practical entropy-compressed rank/select dictionary. *Proceedings of ALENEX*, 2007.
- [8] P. Elias. Efficient storage and retrieval by content and address of static files. *Journal of the ACM*, 1974.
- [9] P. Elias. Universal codeword sets and representations of the integers. *IEEE Transactions on Information Theory*, (21):194–203, 1975.
- [10] R. Fano. The number of bits required to implement an associative memory. *Computer Structures Group, Project MAC, MIT*, 1971.
- [11] Jacobson G. Space-efficient static trees and graphs. *30th annual IEEE symposium on foundations of computer science*, pages 549–554, 1989.
- [12] Blelloch GE and Farzan A. Succinct representations of separable graphs. *Amir A, Parida L (eds) CPM*, page 138–150, 2010.
- [13] Raman R Geary RF, Rahman N. A simple optimal representation for balanced parentheses. *Theor Comput Sci* 368, page 231–246, 2006.
- [14] Simon Gog. *SDSL - Succinct Data Structure Library*. Available at <https://github.com/simongog/sdsl-lite>.
- [15] V. Raman J. I. Munro. Succinct representation of balanced parentheses, static trees and planar graphs. *38th FOCS*, pages 118–126, 1997.

- [16] Daniel K.Blandfor, Guy E. Blelloch, and Ian A. Ksh. Compact representations of separable graphs. *Proc. SODA, ACM/SIAM*, page 679–688, 2003.
- [17] R. J. Lipton, D. J. Rose, and R. E. Tarjan. Generalized nested dissection. *SIAM Journal on Numerical Analysis*, (16):346–358, 1979.
- [18] G. L. Miller, S.-H. Teng, W. P. Thurston, and S. A.Vavasis. Separators for sphere-packings and nearest neighbor graphs. *Journal of the ACM*, (44):1–29, 1997.
- [19] D. Okanohara and K. Sadakane. Practical entropy-compressed rank/select dictionary. *Proceedings of ALENEX*, 2007.
- [20] R. Raman, V. Raman, and S. Rao. Succinct indexable dictionaries with applications to representations of k-ary trees and multi-sets. *SODA*, 2002.
- [21] X.He R.C.N.Chuang, A.Garg and M.Y.Kao. Compact encodings of planar graphs via canonical orderings and multiple parentheses. *Lecture Notes in Computer Science*, 1443, pages 118–129, 1998.
- [22] R.J.Liton and R.E. Tarjan. A separator theorem for planar graphs. *SIAM J.Applied Mathematics*, (16):177–189, 1979.
- [23] svn2github. *Cytoscape GPU layout*. Available at https://github.com/svn2github/cytoscape/tree/master/csplugins/trunk/soc/ghuck/GpuGraphLayoutPlugin/sample_graphs.
- [24] Y.T.Chiang and C.C.Lin. Orderly spanning trees with applications to graph encoding and graph drawing. *In SODA*, pages 506–515, 2001.