

# TCS Technical Report

## A Compact and Fast Index Structure for Families of Sets

by

SHUHEI DENZUMI, JUN KAWAHARA, KOJI TSUDA,  
HIROKI ARIMURA, SHIN-ICHI MINATO,  
AND KUNIHICO SADAKANE

**Division of Computer Science**

**Report Series A**

April 12, 2014

Hokkaido University  
Graduate School of  
Information Science and Technology

Email: [minato@ist.hokudai.ac.jp](mailto:minato@ist.hokudai.ac.jp)

Phone: +81-011-706-7682

Fax: +81-011-706-7682



# A Compact and Fast Index Structure for Families of Sets

Shuhei Denzumi<sup>1</sup>  
Hiroki Arimura<sup>1</sup>

Jun Kawahara<sup>2</sup>  
Shin-ichi Minato<sup>1,4</sup>

Koji Tsuda<sup>3,4</sup>  
Kunihiko Sadakane<sup>5</sup>

<sup>1</sup>Graduate School of IST, Hokkaido University, Japan

<sup>2</sup>Nara Institute of Science and Technology (NAIST), Japan

<sup>3</sup>National Institute of Advanced Industrial Science and Technology (AIST), Japan

<sup>4</sup>ERATO MINATO Discrete Structure Manipulation System Project, JST, Japan

<sup>5</sup>National Institute of Informatics (NII), Japan

{denzumi, arim, minato}@ist.hokudai.ac.jp, jkawahara@is.naist.jp,  
koji.tsuda@aist.go.jp, sada@nii.ac.jp

April 12, 2014

## Abstract

In many real-life problems, we are often faced with manipulating families of sets. Manipulation of large-scale set families is one of the important fundamental techniques for web information retrieval, integration, and mining. For this purpose, a special type of *binary decision diagrams (BDDs)*, called *Zero-suppressed BDDs (ZDDs)*, is used. However, current techniques for storing ZDDs require a huge amount of memory and membership operations are slow. This paper introduces DenseZDD, a compressed index for static ZDDs. Our technique not only indexes set families compactly but also executes fast member membership operations. We also propose a hybrid method of DenseZDD and ordinary ZDDs to allow for dynamic indices.

## 1 Introduction

Binary Decision Diagrams (BDD) [1] are a graph-based representation of Boolean functions and widely used in VLSI logic design and verification. A BDD is constructed reducing a binary decision tree, which represents a decision making process through the input variables. If we fix the order of the input variables and apply the following two reduction rules, then we obtain a minimal and canonical form for a given Boolean function:

1. Delete all redundant nodes (whose two children are identical) and
2. Merge all equivalent nodes (having the same index and pair of children).

Among unique canonical representations of Boolean functions, BDDs are smaller than others such as CNF, DNF, and truth tables for many classes of functions. BDDs have the following features:

- Boolean functions are uniquely represented like other representations.
- Multiple functions are stored compactly by sharing common subgraphs.
- Fast logical operations are executed on Boolean functions.

Zero-suppressed Binary Decision Diagrams (ZDDs) [8] are variation of traditional BDDs, used to manipulate families of sets. Using ZDDs, we can implicitly enumerate combinatorial item set data and efficiently compute set operations over the ZDDs. In the rest of this section, we use the term BDD to indicate both the original BDD and the ZDD unless specified because any ZDD is regarded as a BDD representing some function.

Though BDDs are more compact than other representations of Boolean function and set families, they are still large; a node of a BDD uses 20 to 30 bytes depending on implementations [9]. BDDs become inefficient if the graph size is too large to be held in memory. Therefore the aim of this paper is to reduce the size (number of bits) used to represent BDDs. We classify implementations of BDDs into three types:

- Dynamic: The BDD can be modified. New nodes can be added to the BDD.
- Static: The BDD cannot be modified. Only query operations are supported.
- Freeze-dried: All the information of the BDD is stored, but it cannot be used before restoration.

Most of the current implementations of BDDs are dynamic. There is previous work on freeze-dried representations of BDDs by Starkey and Bryant [16] and later, by Mateu and Prades-Nebot [7]. Hansen, Rao and Tiedemann [4] developed a technique to compress BDD and reduce the size of the BDD to 1-2 bits per node. However there is no implementation of BDDs that is specialized for static case.

This paper is the first to propose a static representation of ZDDs, which we call *DenseZDDs*. The size of ZDDs in our representation is much smaller than an existing dynamic representation [9]. Not only compact, DenseZDD supports much faster membership operations than [9]. Experimental results show that DenseZDDs are five times smaller and membership queries are twenty to several hundred times faster, compared to [9]. Note that our technique can be directly applied to compress traditional BDDs too.

## 2 Preliminaries

Let  $e_1, \dots, e_n$  be items such that  $e_1 < e_2 < \dots < e_n$ . Let  $S = \{a_1, \dots, a_c\}$ ,  $c \geq 0$ , be a set of items. We denote the *size* of  $S$  by  $|S| = c$ . The empty set is denoted by  $\emptyset$ . A *family* is a subset of the power set of all items. A finite family  $F$  of sets is referred to as a *set family*.<sup>1</sup> The *join* of families  $F_1$  and  $F_2$  is defined as  $F_1 \sqcup F_2 = \{S_1 \cup S_2 \mid S_1 \in F_1, S_2 \in F_2\}$ .

<sup>1</sup>In the original ZDD paper by Minato, a set is called a combination, and a set family is called a combinatorial set.

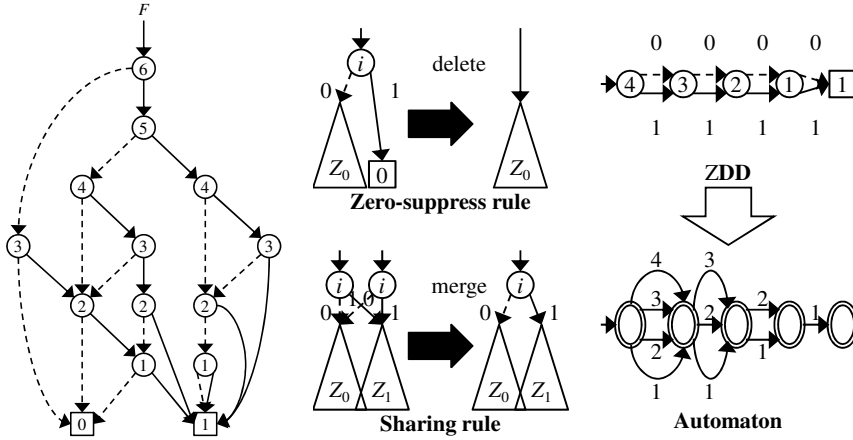


Figure 1: An example of ZDD

Figure 2: Reduction rules of ZDDs

Figure 3: Worst-case example of a straight-forward translation

In the appendix, we describe existing succinct data structures. The balanced parenthesis sequence (BP), the Fully Indexable Dictionary (FID), and some basic structures used in this paper are reviewed. We also explain operations on the data structures such as  $rank_c$ ,  $select_c$ , and so on.

## 2.1 Zero-Suppressed Binary Decision Diagrams (ZDDs)

A *zero-suppressed binary decision diagram* (a ZDD) [8] is a variant of a binary decision diagram [1], customized to manipulate finite families of set. A ZDD is a directed acyclic graph satisfying the following conditions. A ZDD has two types of nodes, terminal and nonterminal nodes. A *terminal node*  $v$  has as attribute a value  $value(v) \in \{0, 1\}$ , indicating whether it is a *0-terminal node* or a *1-terminal node*, denoted by  $\mathbf{0}$  and  $\mathbf{1}$ , respectively. A *nonterminal node*  $v$  has as attributes an integer  $index(v) \in \{1, \dots, n\}$  called the *index*, and two children  $zero(v)$  and  $one(v)$ , called the *0-child* and *1-child*. The edges from nonterminals to their 0-child (1-child resp.) are called *0-edges* (*1-edges* resp.). In the figures, terminal nodes are denoted by squares, and nonterminal nodes are denoted by circles. 0-edges are denoted by dotted arrows, and 1-edges are denoted by solid arrows. We define  $triple(v) = \langle index(v), zero(v), one(v) \rangle$ , called the *attribute triple* of  $v$ . For any nonterminal node  $v$ ,  $index(v)$  is larger than the indices of its children.<sup>2</sup> We define the *size* of the graph, denoted by  $|G|$ , as the number of its nonterminals.

**Definition 1 (set family represented by ZDD)** A ZDD  $G$  rooted at a node  $v \in V$  represents a finite family of sets  $F(v)$  on  $U_n$  defined recursively as follows: (1) If  $v$  is a terminal node:  $F(v) = \{\emptyset\}$  if  $value(v) = 1$ , and  $F(v) = \emptyset$  if  $value(v) = 0$ . (2) If  $v$  is a nonterminal node, then  $F(v)$  is the finite family of sets  $F(v) = (\{e_{index(v)}\} \sqcup F(one(v))) \cup F(zero(v))$ .

<sup>2</sup>In ordinary BDD or ZDD papers, the indices are in ascending order from roots to terminals. For convenience, we employ the opposite ordering in this paper.

Table 1: Main operations supported by ZDD. The first group is the primitive ZDD operations used to implement the others, yet they could have other uses

$index(v)$	Returns the index of node $v$ .
$zero(v)$	Returns the 0-child of node $v$ .
$one(v)$	Returns the 1-child of node $v$ .
$getnode(i, v_0, v_1)$	Generates (or makes a reference to) a node $v$ with index $i$ and two child nodes $v_0 = zero(v)$ and $v_1 = one(v)$ .
$topset(v, i)$	Returns a node with the index $i$ reached by traversing only 0-edges. If such a node does not exist, return the 0-terminal node.
$member(v, S)$	Returns <i>true</i> if $S \in F(v)$ , and returns <i>false</i> otherwise.
$count(v)$	Returns $ F(v) $ .
$offset(v, i)$	Returns $v$ such that $F(v) = \{ S \subseteq U_n \mid S \in F, e_i \notin S \}$ .
$onset(v, i)$	Returns $v$ such that $F(v) = \{ S \setminus \{e_i\} \subseteq U_n \mid S \in F, e_i \in S \}$ .
$apply_{\diamond}(v_1, v_2)$	Returns $v$ such that $F(v) = F(v_1) \diamond F(v_2)$ , for $\diamond \in \{\cup, \cap, \setminus, \oplus\}$ .

The example in Fig. 1 represents a sets family  $F = \{ \{6, 5, 4, 3\}, \{6, 5, 4, 2\}, \{6, 5, 4, 1\}, \{6, 5, 4\}, \{6, 5, 2\}, \{6, 5, 1\}, \{6, 5\}, \{6, 4, 3, 2\}, \{6, 4, 3, 1\}, \{6, 4, 2, 1\}, \{6, 2, 1\}, \{3, 2, 1\} \}$ .

A set  $S = \{a_1, \dots, a_c\}$  describes a path in the graph  $G$  starting from the root. At each nonterminal node, the path continues to the 0-child if  $e_i \notin S$  and to the 1-child if  $e_i \in S$ . The path eventually reaches the 1-terminal (or 0-terminal resp.), indicating that  $S$  is accepted (or rejected resp.).

In ZDD, we employ the following two reduction rules to compress the graph:

- Zero-suppress rule: A nonterminal node whose 1-child is the 0-terminal node.
- Sharing rule: Two or more nonterminal nodes having the same attribute triple. By applying above rules, we can reduce the graph without changing its semantics. If we apply the two reduction rules as much as possible, then we obtain a canonical form for a given family of sets.

We can reduce the size of ZDDs by using a type of *attributed edges* [11] named *0-element edges*. By using the 0-element edge, which simplifies construction of DenseZDD. We can implement 0-element edges by 1-edges with *empty set flags*. This attribute indicates that the pointing subgraph includes the empty set  $\emptyset$  in the family represented by the subgraph. If it is set, it means that the node pointed to by the 1-edge represents a set family including the empty set  $\emptyset$ . We have to place a couple of constraints on using 0-element edges to keep the uniqueness of the graphs: (1) Use the 0-terminal node only. (2) Do not use 0-element edges at the 0-edge on each node. Each nonterminal node  $v$  has an  $\emptyset$ -flag  $empflag(v)$  on its 1-edge to implement 0-element edges. If  $empflag(v) = 1$ , the subgraph pointed by the  $v$ 's 1-edge includes the empty set  $\emptyset$  in the family represented by the subgraph. In this paper, effective  $\emptyset$ -flags are denoted as circles at starting points of 1-edges.

Table 1 summarizes operations of ZDDs. The upper half shows the primitive operations, while the lower half shows other operations which can be implemented by using the primitive operations. The operations  $index(v)$ ,  $zero(v)$ ,  $one(v)$ ,  $topset(v, i)$  and  $member(v, S)$  do not create new nodes. Therefore they can be done on a static ZDD. The operation  $count(v)$  does not create any node; however we need an auxiliary array to memorize which nodes are already visited.



**Theorem 1** *Let  $u$  be the size of the ZDD that removes the zero-suppress rule only for nodes pointed to by 0-edges. A ZDD with  $m$  nodes on  $n$  items can be stored in  $2u + m \log m + 3m + o(u)$  bits so that the primitive operations except  $\text{getnode}(i, v_0, v_1)$  are done in constant time. In other words,  $u$  is the size of the ZDD with dummy nodes that are described below. The  $\text{getnode}(i, v_0, v_1)$  operation is done in  $\mathcal{O}(\log m)$  time.*

**Theorem 2** *A ZDD with  $m$  nodes on  $n$  items can be stored in  $\mathcal{O}(m(\log m + \log n))$  bits so that the primitive operations are done in  $\mathcal{O}(\log m)$  time except  $\text{getnode}(i, v_0, v_1)$ . The  $\text{getnode}(i, v_0, v_1)$  operation is done in  $\mathcal{O}(\log^2 m)$  time.*

### 3.1 DenseZDD

A DenseZDD  $DZ = \langle U, M, I \rangle$  is composed of three data structures: a zero-edge tree  $U$ , a dummy node vector  $M$ , and a one-child array  $I$ .

**Zero-edge tree** : The spanning tree of ZDD  $G$  formed by the 0-edges is called the *zero-edge tree* of  $G$  and denoted by  $T_Z$ . In a zero-edge tree, all 0-edges are reversed and the 0-terminal node becomes the root of the tree. The preorder rank of each node is used to identify it. Zero-edge trees are based on the same idea as left or right trees by Maruyama *et al.* [6].

An important difference between our structure and theirs is the existence of *dummy nodes*. We call nodes in the original ZDD as *real nodes*. We use the zero-edge tree with dummy nodes, denoted by  $T'_Z$ . We create dummy nodes on each 0-edge to guarantee that the depth of every real node  $v$  in the zero-edge tree equals  $\text{index}(v)$ . We define the depth of the 0-terminal node, the root of this tree, to be 0. Let  $U$  be the BP of  $T'_Z$ . The length of  $U$  is  $\mathcal{O}(mn)$  because we create  $n - 1$  dummy nodes for one real node in the worst case. An example of a zero-edge tree and its BP are shown in Fig. 5. Black circles are dummy nodes and the number next to each node is its preorder rank. The 0-terminal node is ignored in the BP because we know the root of a zero-edge tree is always that node.

**Dummy node vector** : A bit vector of the same length as  $U$  is used to distinguish dummy nodes and real nodes. We call it the *dummy node vector* of  $T'_Z$  and denote it by  $B_D$ . The  $i$ -th bit is 1 if and only if the  $i$ -th parenthesis of  $U$  is '(' and its corresponding node is a real node in  $T'_Z$ . An example of a dummy node vector is also shown in Fig. 5. The 0-terminal node is also ignored. Let the FID of  $B_D$  be  $M$ . Using  $M$ , we can determine whether a node is dummy or real, and compute preorder ranks among only real nodes. We can also obtain positions of real nodes on BP from their preorder ranks by the select operation on  $M$ .

**One-child array** : An integer array to represent the 1-child of each node is called the *one-child array* and denoted by  $C_O$ . This array contains node preorder ranks of all 1-children in preorder on  $T_Z$ . That is, its  $i$ -th element is the preorder rank of the 1-child of the nonterminal node whose preorder rank is  $i$ . We also require one bit for each element of the one-child array to store the  $\emptyset$ -flag. If  $\text{empflag}(v) = 1$  for a nonterminal node  $v$ , the corresponding element in the one-child array will be negative. An example of a one-child array is shown



in Fig. 6. Let  $I$  be the compressed representation of  $C_O$ . In  $I$ , one integer is represented by  $\lceil \log(m+1) \rceil + 1$  bits, including one bit for the  $\emptyset$ -flag.

## 4 Algorithm

### 4.1 Conversion of an Ordinary ZDD to a DenseZDD

We show how to construct the DenseZDD. We first build the zero-edge tree from the given ZDD. A pseudo-code is given in Fig. 9 in the appendix. The zero-edge tree consists of all 0-edges of the ZDD, with their directions being reversed. For a nonterminal node  $v$ , we say that  $v$  is a  $0^r$ -child of  $zero(v)$ . To make a zero-edge, we use a list *revzero* in each node, which stores  $0^r$ -children of the node. The lists for all the nodes are computed by a depth-first traversal of the ZDD. This is done in  $\mathcal{O}(m)$  time and  $\mathcal{O}(m)$  space, since each node is visited at most twice and the total size of *revzero* is the same as the number of nonterminal nodes.

We obtained a zero-edge tree  $T$ , but it is not an ordered tree. We define preorder rank  $prank(v)$  for every node  $v$  before traversal. The nodes in *revzero* are sorted in descending order of their pairs  $\langle index, prank \rangle$ , that is,  $index(revzero[i]) \geq index(revzero[i+1])$  for  $1 \leq i < |revzero(v)|$ . Then, nodes with higher indices are visited first. This ordering is useful to reduce the number of dummy nodes and to implement ZDD operations simply. It seems impossible to define visiting order of nodes by preorder rank of their 1-children during computing preorder, but it is possible. Since a ZDD node  $v$  satisfies  $index(v) > index(zero(v))$  and  $index(v) > index(one(v))$ , we can decide  $prank$  for every node by the pseudo code in Fig. 7, which is a BFS algorithm based on *index* value starting from 0-terminal. To compute  $prank$  efficiently, we construct the temporary BP for the zero-edge tree. Using the BP, we can compute the size of each subtree rooted by  $v$  in  $T$  in constant time and compact space.

Next, we create dummy nodes imaginarily. For a node  $v$ , we create  $q = \max\{i \in \{1, \dots, n\} \mid i = index(revzero[j]) - 1, 1 \leq j \leq |revzero(v)|\}$  dummy nodes  $d_1, \dots, d_q$  such that  $triple(d_i) = \langle index(v) + i, d_{i-1}, \mathbf{0} \rangle$ , and  $empflag(d_i) = 0, 1 \leq i \leq q$ . For convention,  $d_0$  denotes  $v$ .

To sum up, the DenseZDD for the given ZDD is composed of the zero-edge tree, the one-child array, and the dummy node vector. We traverse the zero-edge tree in DFS order as if dummy nodes exist and construct the BP representation  $U$ , the dummy node vector  $M$ , and the one-child array  $I$ . The BP and dummy node vector are constructed for the zero-edge tree with dummy nodes. On the other hand, the one-child array ignores dummy nodes. DenseZDD  $DZ = \langle U, M, I \rangle$  is obtained. Pseudo-codes are given in algorithms in Fig. 8, and 9 in the appendix.

### 4.2 Primitive ZDD Operations

We show how to implement primitive ZDD operations on DenseZDD  $DZ = \langle U, M, I \rangle$  except *getnode*. We give an algorithm for *getnode* in Section 5.

In the zero-edge tree, there are two types of nodes: real nodes and dummy nodes. Real nodes are those in the ZDD, while dummy nodes have no corresponding ZDD nodes. Real nodes are numbered from 1 to  $m$  based on preorders

in the tree. Below a node is identified with this number, which we call its *node number*. We can convert between the node number  $i$  of a node and the position  $p$  in the BP sequence  $U$  by  $p := \text{select}_1(M, i)$  and  $i := \text{rank}_1(M, p)$ . The 0-terminal has node number 0 and nonterminal nodes have positive node numbers. If a node number of a negative value is used, it means a node with an  $\emptyset$ -flag.

In addition, we consider an additional primitive operation for DenseZDDs:  $\text{chkdum}(p)$ . This operation checks if a node at position  $p$  on  $U$  is a dummy node or not. If it is a dummy  $\text{chkdum}$  returns false; otherwise it returns true. This operation is implemented by simply looking at the  $p$ -th bit of  $M$ . If the bit is 0, then the node is dummy; otherwise it is a real node.

$\text{index}(i)$  : Since the item of the node is the same as the depth of the node, we can obtain  $\text{index}(i) := \text{depth}(U, \text{select}_1(M, i))$ .

$\text{one}(i)$  :  $\text{one}(i)$ : Because 1-children are stored in preorder of the parents of nodes, we can obtain  $\text{one}(i) := I[i]$ .

$\text{topset}(i, d)$  : The node  $\text{topset}(i, d)$  is the ancestor of node  $i$  in the zero-edge tree with index  $d$ . A naive solution is to iteratively climb up the zero-edge tree from node  $i$  until we reach a node with index  $d$ . However, as shown above, the index of a node is identical to its depth. By using the power of the succinct tree data structure, we can directly find the answer by  $\text{topset}(i, d) := \text{rank}_1(M, \text{level\_ancestor}(U, \text{select}_1(M, i), d))$ .

$\text{zero}(i)$  : Implementing the  $\text{zero}$  operation requires a more complicated technique. Consider a subtree  $T$  of the zero-edge tree consisting of the node  $i$ , its real parent node  $r$ , all real children of  $r$ , and dummy nodes between those nodes. As a pre-condition, the zero-edge tree is constructed by Algorithm 9 in the appendix. That is, for the children of  $r$ , the nodes with higher  $\text{index}$  value have smaller preorder, and the imaginary parents of the children are dummy nodes (or  $i$ ) that are added on the edge between  $r$  and the child having the highest  $\text{index}$  value. Computing  $\text{zero}(i)$  is equivalent to finding  $r$ . Because the children of  $r$  are ordered from left to right in descending order of their depths, and dummy nodes are shared as much as possible, the deepest node in  $T$  is on the leftmost path from  $r$ . Furthermore, the parents of other real children are also on the leftmost path. This property also holds in the original zero-edge tree. The dummy node vector  $B_D$  stores flags in the preorder in the zero-edge tree. Then  $B_D[p_r] = B_D[p_i] = 1$ , where  $p_r$  and  $p_i$  are positions of nodes  $r$  and  $i$  in the BP sequence  $U$ , and  $B_D[j] = 0$  for any  $p_r < j < p_i$ . Therefore we can find  $p_r$  by a rank operation. In summary,  $\text{zero}(i) := \text{rank}_1(M, \text{parent}(U, \text{select}_1(M, i)))$ .

### 4.3 Compressing the Balanced Parentheses Sequence

The balanced parentheses sequence  $U$  is of length  $2u$ , where  $u$  is the number of nodes including dummy nodes. Let a ZDD have  $m$  real nodes and the number of items be  $n$ ,  $u$  is  $mn$  in the worst case. Here we compress the BP sequence  $U$ .

The BP sequence  $U$  consists of at most  $2m$  runs of identical symbols. To see this, consider the substring of  $U$  between the positions for two real nodes. There

is a run ‘))...’ followed by a run ‘(((...’ in the substring. We encode lengths of those runs using some integer encoding scheme such as the delta-code or the gamma-code [3]. An integer  $x > 0$  is encoded in  $\mathcal{O}(\log x)$  bits. Because the maximum length of a run is  $n$ ,  $U$  can be encoded in  $\mathcal{O}(m \log n)$  bits. The range min-max tree of  $U$  has  $2m/\log m$  leaves. Each leaf corresponds to a substring of  $U$  that contains  $\log m$  runs. Then any tree operation can be done in  $\mathcal{O}(\log m)$  time. The range min-max tree is stored in  $\mathcal{O}(m(\log n + \log m)/\log m)$  bits.

We also compress the dummy node vector  $B_D$ . Because its length is  $2u \leq 2mn$  and there are only  $m$  ones, it can be compressed in  $m(2 + \log m) + o(u)$  bits by FID. The operations  $select_1$  and  $rank_1$  take constant time. We can reduce the term  $o(u)$  to  $o(m)$  by using a sparse array [14]. The operation  $select_1$  is done in constant time, while  $rank_1$  takes  $\mathcal{O}(\log m)$  time. From the discussions above, we can prove Theorem 1 and Theorem 2. For the proof, see the appendix.

## 5 Hybrid method

In this section, we show how to implement dynamic operations on DenseZDD. Namely, we need to implement the  $getnode(i, v_0, v_1)$  operation. Our approach is to use a hybrid data structure using both the DenseZDD and a conventional dynamic ZDD. Assume that initially all the nodes are in a DenseZDD. Let  $m_0$  be the number of initial nodes. In a dynamic ZDD, the operation  $getnode(i, v_0, v_1)$  is implemented by a hash table indexed with the triple  $\langle i, v_0, v_1 \rangle$ .

We show first how to check whether the node  $v := getnode(i, v_0, v_1)$  already exists. That is, we want to find a node  $v$  such that  $index(v) = i$ ,  $zero(v) = v_0$ ,  $one(v) = v_1$ . If  $v$  does not exist, we create such a node using the hash table as well as a dynamic ZDD. If it exists, in the zero-edge tree,  $v$  is a real child node of  $v_0$ . Consider again the subtree of the zero-edge tree rooted at  $v_0$  and having all real children of  $v_0$ . All children of  $v_0$  with index  $i$  share the common (possible dummy) parent node, say  $w$ . Because  $w$  is on the leftmost path in the subtree, it is easy to find it. Namely,  $w := level\_ancestor(U, select_1(M, rank_1(M, v_0) + 1), i)$ . The node  $v$  is a child of  $w$  with  $one(v) = v_1$ . Because all children of  $w$  are sorted in the order of  $one$  values by the construction algorithms, we can find  $v$  by a binary search. For this, we use  $degree$  and  $child$  operations on the zero-edge tree.

**Theorem 3** *The existence of  $getnode(i, v_0, v_1)$  can be checked in  $\mathcal{O}(t \log m)$  time, where  $t$  is the time complexity of primitive ZDD operations.*

If the BP sequence is not compressed,  $getnode$  takes  $\mathcal{O}(\log m)$  time. Otherwise it takes  $\mathcal{O}(\log^2 m)$  time. We should check the hash table before checking the zero-edge tree if dynamic nodes are already exist. As well as a conventional ZDD, hashing increases constant factors of time bounds significantly and add space bound  $\mathcal{O}(x \log x)$  where  $x$  is the number of dynamic nodes.

## 6 Experimental Results

We ran experiments to evaluate the compression, construction, and operation times of DenseZDDs. We implemented the algorithms described in Sec. 3 and 4

Table 2: Comparison of performance, where  $\delta$  denotes the dummy node ratio

data set	#items	#nodes	#itemsets	size (bytes)			comp. ratio		$\delta$
				Z	DZ	DZ <sub>c</sub>	DZ	DZ <sub>c</sub>	
<i>grid5</i>	40	584	8,512	17,520	2,350	2,196	0.134	0.126	0.28
<i>grid10</i>	180	377,107	$4.1 \times 10^{20}$	11,313,210	1,347,941	1,265,773	0.119	0.112	0.20
<i>grid15</i>	420	$1.5 \times 10^8$	$2.3 \times 10^{48}$	4,342,789,110	678,164,945	647,843,001	0.156	0.149	0.19
<i>webview5</i>	952	2,299	11,928	10,592,760	3,871,679	1,851,889	0.365	0.174	0.93
<i>webview10</i>	1,617	6,060	70,713	281,700	1,034,471	477,299	0.367	0.169	0.93
<i>webview20</i>	2,454	30,413	634,065	912,390	290,661	140,873	0.318	0.154	0.92
<i>webview50</i>	2,905	93,900	$4.4 \times 10^6$	181,800	44,846	24,596	0.246	0.135	0.88
<i>webview100</i>	3,149	353,092	$2.7 \times 10^7$	68,970	11,455	7,967	0.166	0.115	0.75
<i>webviewALL</i>	3,149	465,449	$3.2 \times 10^7$	13,963,470	4,964,303	2,413,625	0.355	0.172	0.92
<i>randjoin128</i>	32,696	6,751	$2.5 \times 10^8$	202,530	408,149	99,117	2.015	0.489	0.99
<i>randjoin2048</i>	32,768	377,492	$1.8 \times 10^{13}$	11,324,760	2,415,648	1,511,658	0.213	0.133	0.82
<i>randjoin8192</i>	32,768	$1.3 \times 10^6$	$3.7 \times 10^{15}$	38,094,930	5,328,502	4,386,452	0.139	0.115	0.42
<i>randjoin16384</i>	32,768	$1.9 \times 10^6$	$2.8 \times 10^{16}$	56,447,280	7,056,418	6,113,910	0.125	0.108	0.14

Table 3: Converting time and random searching time

data set	conversion time (sec)				traverse time (sec)			search time (sec)		
	read	convert	const.	comp.	Z	DZ	DZ <sub>c</sub>	Z	DZ	DZ <sub>c</sub>
<i>grid5</i>	0.001	0.001	0.009	0.000	0.000	0.001	0.006	0.029	0.038	0.229
<i>grid10</i>	0.461	0.634	0.449	0.060	0.075	0.247	1.388	0.005	0.013	0.056
<i>grid15</i>	124.887	407.502	112.379	8.186	41.214	102.673	398.397	0.006	0.011	0.064
<i>webview5</i>	0.256	0.690	1.361	0.055	0.066	0.154	0.250	1.966	0.045	0.099
<i>webview10</i>	0.217	0.226	0.564	0.041	0.017	0.042	0.073	1.901	0.043	0.100
<i>webview20</i>	0.066	0.036	0.313	0.022	0.005	0.014	0.027	1.875	0.046	0.101
<i>webview50</i>	0.013	0.019	0.050	0.004	0.002	0.020	0.005	1.314	0.273	0.102
<i>webview100</i>	0.004	0.002	0.017	0.001	0.000	0.004	0.007	0.777	0.129	0.376
<i>webviewALL</i>	0.551	0.927	1.644	0.108	0.091	0.207	0.346	1.706	0.049	0.105
<i>randjoin128</i>	0.004	0.053	0.149	0.008	0.001	0.002	0.003	0.527	0.044	0.095
<i>randjoin2048</i>	0.243	0.742	0.946	0.029	0.093	0.126	0.145	8.071	0.044	0.098
<i>randjoin8192</i>	0.858	2.573	1.259	0.043	0.338	0.240	0.304	15.604	0.039	0.092
<i>randjoin16384</i>	1.270	5.016	1.471	0.070	0.676	0.353	0.447	19.501	0.040	0.093

in C/C++ languages on top of the SAPPORO BDD package [9]. The package is general implimentation of ZDD with 0-element edges, and uses 30 bytes per ZDD node. We performed experiments on eight quad-core 3.09 GHz AMD Opteron 8393 SE processors (i.e, 32 CPU cores in total) and 512 GB DDR2 memory shared among cores running. SUSE 10 Our algorithms use a single core since they are not parallelized.

We show the characteristics of the ZDDs in Table 2. Original ZDDs are denoted by Z. DenseZDDs without/with compression of the balanced parentheses sequences of the zero-edge trees are denoted by DZ/DZ<sub>c</sub>, respectively.

As real data sets, for  $N = 5, 10, 20, 50, 100$ , the source ZDD *webviewN* was constructed from the data set BMS-Web-View-2<sup>3</sup> by using mining algorithm LCM over ZDD [12] with minimum support  $N$ . For artificial data sets, the ZDD *gridN* represents all self-avoiding paths on an  $N \times N$  grid graph from the top left corner to the bottom right corner [5]<sup>4</sup>. Finally, *randjoinN* is a ZDD that represents the join  $C_1 \sqcup \dots \sqcup C_4$  of four ZDDs for random families  $C_1, \dots, C_4$

<sup>3</sup><http://fimi.ua.ac.be>

<sup>4</sup>An algorithm animation: <http://www.youtube.com/watch?v=Q4gTV4r0zRs>

consisting of  $N$  sets of size one drawn from the set of  $n = 32768$  items.

In Table 2 we show the sizes of the original ZDD, the DenseZDD with/without compression and their compression ratio. We compressed FID for dummy node vector if the dummy node ratio is more than 75%. In almost cases, we observe that the DenseZDD is from 2.5 to 9 times smaller than the original ZDD, and that compressed DenseZDD is from 6 to 9 times smaller than the original ZDD. The compressed DenseZDD is quarter the size of the DenseZDD in the best case, and is half the size of the original ZDD in the worst case. For most of our data sets, the ratio  $\delta$  of the number of dummy nodes to the size of DenseZDD is roughly 90%, except for *gridN* and *randjoin16384*.

In Table 3, we show the conversion times from ZDDs to DenseZDDs, traversal times, and search times on ZDDs and DenseZDDs. Conversion time is composed of four parts: time to read a file containing a stored ZDD and reconstruct the ZDD, convert it to raw parentheses, bits, and integers, construct succinct representation of them, and compress the BP of the zero-edge tree. The conversion time appears almost linear in the input size showing its scalability for large data. Traverse operation used *zero(v)* and *one(v)*, while membership operation used *topset(v, i)* and *one(v)*. We observed that the DenseZDD has almost twice longer traverse time and more than 10 times shorter search time than an original ZDD. These results show the efficiency of our implementation of the *topset(v, i)* operation on DenseZDD using level-ancestor operations.

From the above results, we conclude that DenseZDDs are more compact than ordinary ZDDs unless the dummy node ratio is extremely high, and the membership operations for DenseZDDs are much faster if the number of items is large or the dummy node ratio is small. We observed that in DenseZDDs, traversal time is approximately double and search time approximately one-tenth compared to the original ZDDs. The traversal is accelerated especially for large-scale sparse datasets because the number of nodes connected by 0-edges grows as large as the index number  $n$ . Recently, processing of “Big Data” have attracted a great deal of attention, and we often deal with a large-scale sparse dataset, which has more than ten thousands of items as the columns of a dataset. In the era of Big Data, we expect that DenseZDD will be effective for various real-life applications, such as data mining, system diagnosis, and network analysis.

## 7 Conclusion

In this paper, we have presented a compressed index for static ZDDs named DenseZDD. We also proposed a hybrid method for dynamic operations on DenseZDD so that we can manipulate DenseZDD and conventional ZDD together. For future work, the one-child array should be stored in more compact space. Constructing the DenseZDD from the normal ZDD using external memory is an important open problem. We will implement the hybrid method on our ZDD package and convert/update algorithm with less memory. We expect that our technique can be extended to other variants of BDDs.

## Acknowledgments

The authors would like to thank Prof. Roberto Grossi, Prof. Rajeev Raman, and Dr. Yasuo Tabei for their discussions and valuable comments. This work was supported by Grant-in-Aid for JSPS Fellows 25-193700. This research was

partly supported by Grant-in-Aid for Scientific Research on Innovative Areas — Exploring the Limits of Computation, MEXT, Japan and ERATO MINATO Discrete Structure Manipulation System Project, JST, Japan. KT is supported by JST CREST and JSPS Kakenhi 25106005. KS is supported by JSPS KAKENHI 23240002.

## References

- [1] Randal E. Bryant. Graph-based algorithms for Boolean function manipulation. *IEEE Transactions on Computers*, C-35(8):677–691, 1986.
- [2] Shuhei Denzumi, Ryo Yoshinaka, Hiroki Arimura, and Shin-ichi Minato. Notes on sequence binary decision diagrams: Relationship to acyclic automata and complexities of binary set operations. In *Prague Stringology Conference 2011*, pages 147–161, Prague, 2011.
- [3] Peter Elias. Universal codeword sets and representation of the integers. *IEEE Transactions on Information Theory*, IT-21(2):194–203, 1975.
- [4] Esben Rune Hansen, S. Srinivasa Rao, and Peter Tiedemann. Compressing binary decision diagrams. In *18th European Conference on Artificial Intelligence*, pages 799–800, 2008.
- [5] Donald E. Knuth. *Combinatorial Algorithms, part 1: Volume 4A of The Art of Computer Programming*. Addison-Wesley Professional, Boston, 1st edition, 2011.
- [6] Shirou Maruyama, Masaya Nakahara, Naoya Kishiue, and Hiroshi Sakamoto. ESP-index: A compressed index based on edit-sensitive parsing. *Journal of Discrete Algorithms*, 2013.
- [7] P. Mateu-Villarroya and J. Prades-Nebot. Lossless image compression using ordered binary-decision diagrams. *Electronics Letters*, 37:162–163, 2001.
- [8] Shin-ichi Minato. Zero-suppressed BDDs for set manipulation in combinatorial problems. In *30th international Design Automation Conference*, pages 272–277, 1993.
- [9] Shin-ichi Minato. SAPPORO BDD package. Division of Computer Science, Hokkaido University, 2012. to be released.
- [10] Shin-ichi Minato. Z-skip-links for fast traversal of zdds representing large-scale sparse datasets. In *ESA 2013*, volume 8125 of *LNCS*, pages 731–742, Heidelberg, 2013. Springer.
- [11] Shin-ichi Minato, Nagisa Ishiura, and Shuzo Yajima. Shared binary decision diagram with attributed edges for efficient Boolean function manipulation. In *27th international Design Automation Conference*, pages 52–57, 1990.
- [12] Shin-ichi Minato, Takeaki Uno, and Hiroki Arimura. LCM over ZBDDs: Fast generation of very large-scale frequent itemsets using a compact graph-based representation. In *Advances in Knowledge Discovery and Data Mining*, volume 5012 of *LNCS*, pages 234–246, Heidelberg, 2008. Springer.

- [13] Gonzalo Navarro and Kunihiro Sadakane. Fully-functional static and dynamic succinct trees. *ACM Transactions on Algorithms*, 2010. Accepted A preliminary version appeared in SODA 2010.
- [14] Daisuke Okanohara and Kunihiro Sadakane. Practical entropy-compressed rank/select dictionary. In *Ninth Workshop on Algorithm Engineering and Experiments*, pages 60–70, 2007.
- [15] Rajeev Raman, Venkatesh Raman, and Srinivasa Rao Satti. Succinct indexable dictionaries with applications to encoding  $k$ -ary trees, prefix sums and multisets. *ACM TALG*, 3(4):43:1–43:25, 2007.
- [16] Mike Starkey and Randy Bryant. Using ordered binary-decision diagrams for compressing images and image sequences. Technical Report CMU-CS-95-105, Carnegie Mellon University, 1995.

## A Succinct data structures

### A.1 Succinct data structures for rank/select

Let  $B$  be a binary vector of length  $u$ , that is,  $B[i] \in \{0, 1\}$  for any  $0 \leq i < u$ . The rank value  $\text{rank}_c(B, i)$  is defined as the number of  $c$ 's in  $B[0..i]$ , and the select value  $\text{select}_c(B, j)$  is the position of  $j$ -th  $c$  ( $j \geq 1$ ) in  $B$  from the left. Note that  $\text{rank}_c(B, \text{select}_c(B, j)) = j$  holds if  $j \leq \text{rank}_c(B, n - 1)$ , the number of  $c$ 's in  $B$ . The predecessor  $\text{pred}_c(B, i)$  is defined as the position  $j$  of the rightmost  $c = B[j]$  to the left of  $B[i]$ . The predecessor is computed by  $\text{pred}_c(B, i) := \text{select}_c(B, \text{rank}_c(B, i))$ .

The Fully Indexable Dictionary (FID) is a data structure for computing rank and select on binary vectors [15].

**Theorem 4 (Raman et al. [15])** *For a binary vector of length  $u$  with  $n$  ones, its Fully Indexable Dictionary uses  $\binom{u}{n} + \mathcal{O}(u \log \log u / \log u)$  bits of space and computes  $\text{rank}_c(B, i)$  and  $\text{select}_c(B, i)$  in constant time on the  $\Omega(\log u)$ -bit word RAM.*

This data structure uses asymptotically optimal space because any data structure for storing the vector uses  $\lceil \binom{u}{n} \rceil$  bits in the worst case. Such a data structure is called a *succinct data structure*.

### A.2 Succinct data structures for trees

An ordered tree is a rooted unlabeled tree such that children of each node have some order. A succinct data structure for an ordered tree with  $n$  nodes uses  $2n + o(n)$  bits of space and supports various operations on the tree such as finding the parent or  $i$ -th child, computing the depth or the preorder of a node, etc., in constant time [13]. An ordered tree with  $n$  nodes is represented by a string of length  $2n$  called a balanced parentheses sequence (BP), defined by a depth-first traversal of the tree. Starting from the root, we write an open parenthesis '(' if we arrive at a node from above, and a close parenthesis ')' if we leave from a node upward.

In this paper, the following operations are used. Let  $P$  denote the BP sequence of a tree. A node is identified with the position of the open parenthesis in  $P$  representing the node.

- $\text{depth}(P, i)$ : the depth of a node at position  $i$ . (The depth of a root is 0.)
- $\text{preorder}(P, i)$ : the preorder of a node at position  $i$ .
- $\text{level\_ancestor}(P, i, d)$ : the position of the ancestor with depth  $d$  of node  $i$ .
- $\text{parent}(P, i)$ : the position of the parent of node  $i$  (identical to  $\text{level\_ancestor}(P, i, \text{depth}(P, i) - 1)$ ).
- $\text{degree}(P, i)$ : the number of children of node  $i$ .
- $\text{child}(P, i, d)$ : the  $d$ -th child of node  $i$ .



The operations take constant time.

A brief overview of the data structure is the following. The BP sequence is partitioned into equal-length blocks. The blocks are stored in leaves of a rooted tree called range min-max tree. In each leaf of the range min-max tree, we store the maximum and the minimum values of node depths in the corresponding block. In each internal node, we store the maximum and the minimum of values stored in children of the node. By using this range min-max tree, all tree operations are implemented efficiently.

## B Proof

Following is the proof for Theorem 1 and 2.

**Proof 1** *We first prove Theorem. 1. From the above discussion, the BP  $U$  of zero-edge tree costs  $2u = \mathcal{O}(mn)$  bits where  $u$  is the size of corresponding quasi-ZDD. The one-child array needs  $m \log m$  bits for 1-children and  $m$  bits for  $\emptyset$ -flags. Using FID, the dummy node vector is stored in  $m(2 + \log m) + o(u)$  bits. Therefore, the DenseZDD can be stored in  $2u + m \log m + 3m + o(u)$  bits and primitive operations except `getnode` are done in constant time because the  $\text{rank}_1$ ,  $\text{select}_1$ , and any tree operations take constant time. Since the `getnode` finds a target node by binary search, it takes  $\mathcal{O}(\log m)$  (described in Sec. 5).*

*Next, we prove Theorem 2. When we compress  $U$ , it can be stored in  $\mathcal{O}(m \log n)$  bits and the min-max tree is stored in  $\mathcal{O}(m(\log n + \log m)/\log m)$  bits. The dummy node vector can be compressed in  $m(2 + \log m) + o(m)$  bits by FID with sparse array. But, the time order of any tree operations and the  $\text{rank}_1$  operation is changed from constant time to  $\mathcal{O}(\log m)$  time. Therefore, the DenseZDD can be stored in  $\mathcal{O}(m(\log m + \log n))$  bits and primitive operations take  $\mathcal{O}(\log m)$  times larger than the above time because all of them use tree operations or  $\text{rank}_1$  on  $M$ .*

## C Pseudo codes

### C.1 Algorithm to Compute Preorder of Nodes in Advance

**Global variables:**  $L_1, \dots, L_n$  are list which are empty initially.

**ALGORITHM Compute\_Preorder** ( $L_0$ )

**Input:**  $L_0$ : a list stores only  $\langle \{\mathbf{0}\}, [0, stsize(\mathbf{0}) - 1] \rangle$ ;

- 1: **for**  $i = 0, \dots, n$
- 2:   **for each**  $\langle A, [l, r] \rangle \in L_i$  in arbitrary order %  $A$  is a set of nodes
- 3:     **for each**  $v \in A$  in descending order of  $\langle prank(one(v)), empflag(v) \rangle$
- 4:        $prank(v) \leftarrow l++$ ;
- 5:       **for each**  $j \in \{j \mid w \in revzero(v), j = index(w)\}$  in descending order
- 6:          $A \leftarrow \{w \mid w \in revzero(v), index(w) = j\}$ ;
- 7:          $r \leftarrow l + \text{sum}\{stsize(w) \mid w \in B\}$ ;
- 8:         append  $\langle B, [l, r] \rangle$  to  $L_j$ ;
- 9:         % That is, the  $prank$  of descendants of nodes in  $B$  are in  $[l, r]$ .
- 9:          $l \leftarrow r + 1$ ;
- 10: **return**;

Figure 7: An algorithm which computes the preorder rank  $prank(v)$  of each node  $v$ . Sets of nodes are implemented by arrays or lists in this code. The  $prank(\mathbf{0})$  is 0.

## C.2 Algorithm to Compute a Zero-edge Tree and a Dummy Node Vector

**ALGORITHM** *Convert\_ZDD\_BitVectors* ( $v, paren, dummy, onechild$ )

**Input:** ZDD node  $v$ , list of parentheses  $paren$ ,  
list of bits  $dummy$ , list of integers  $onechild$

```

1:  $i = index(v)$ ;
2: for each  $w \in revzero(v)$  in ascending order of  $prank(w)$ ;
3:   while  $i + 1 < index(w)$ 
4:     append '(' to  $paren$ , and '0' to  $dummy$ ;
5:      $++i$ ;
6:     append '(' to  $paren$ , and '1' to  $dummy$ ;
7:     append  $prank(one(w)) \cdot (-1 \cdot empflag(w))$  to  $onechild$ ;
8:     Convert_ZDD_BitVectors( $w, paren, dummy, onechild$ );
9:     append ')' to  $paren$ , and '0' to  $dummy$ ;
10: while  $i > index(v)$ 
11:   append ')' to  $paren$ , and '0' to  $dummy$ ;
12:    $--i$ ;
13: return;

```

Figure 8: Algorithm for obtaining the BP representation of the zero-edge tree, the dummy node vector, and the one-child array.

### C.3 Algorithm to Convert a DenseZDD from an Ordinary ZDD

**ALGORITHM** `Construct_DenseZDD` ( $W$ : list of ZDD root)

**Output:** DenseZDD  $DZ$

- 1: **for each**  $v \in V$  compute *revzero* fields for all descendants of  $v$ ;
- 2: compute *stsize* field for all  $0^r$ -descendants;
- 3: **Compute\_Preorder**( $\{\{0\}, [0, \text{stsize}(0) - 1]\}$ );
- 4: create empty lists *paren*, *dummy*, *onechild*;
- 5: append '(' to *paren*, and '0' to *dummy*;
- 6: **Convert\_ZDD\_BitVectors**( $0$ , *paren*, *dummy*, *onechild*);
- 7: append ')' to *paren*, and '0' to *dummy*;
- 8: make BP  $U$  from *paren*;
- 9: make FID  $M$  from *dummy*;
- 10: make compressed representation  $I$  of *onechild*;
- 11: **return**  $DZ \leftarrow \langle U, M, I \rangle$ ;

Figure 9: Algorithm for constructing the DenseZDD from a source ZDD.