

# Compressing Forwarding Tables for Datacenter Scalability

Ori Rottenstreich\*, Marat Radan\*, Yuval Cassuto\*, Isaac Keslassy\*,  
Carmi Arad<sup>‡</sup>, Tal Mizrahi<sup>‡\*</sup>, Yoram Revah<sup>‡</sup> and Avinatan Hassidim<sup>§</sup>

\*Technion, {or@tx, radan@tx, ycassuto@ee, isaac@ee}.technion.ac.il

<sup>‡</sup>Marvell Israel, {carmi, talmi, yoramr}@marvell.com

<sup>§</sup>Google Israel and Bar-Ilan University, avinatanh@google.com

**Abstract**—With the rise of datacenter virtualization, the number of entries in the forwarding tables of datacenter switches is expected to scale from several thousands to several millions. Unfortunately, such forwarding table sizes would not fit on-chip memory using current implementations.

In this paper, we investigate the compressibility of forwarding tables. We first introduce a novel forwarding table architecture with separate encoding in each column. It is designed to keep supporting fast random accesses and fixed-width memory words. Then, we show that although finding the optimal encoding is NP-hard, we can suggest an encoding whose memory requirement per row entry is guaranteed to be within a small additive constant of the optimum. Next, we analyze the common case of two-column forwarding tables, and show that such tables can be presented as bipartite graphs. We deduce graph-theoretical bounds on the encoding size. We also introduce an algorithm for optimal conditional encoding of the second column given an encoding of the first one. In addition, we explain how our architecture can handle table updates. Last, we evaluate our suggested encoding techniques on synthetic forwarding tables as well as on real-life tables.

**Index Terms**—Datacenter Virtualization, Layer-2 Datacenter, Forwarding Information Base, Compression.

## I. INTRODUCTION

### A. Background

The rapid growth of forwarding tables in network switches raises serious concerns in the networking industry. Layer 2 (L2) networks are no longer constrained to small local area networks. On the contrary, they are now used in datacenter networks, which will soon have a need for millions of hosts in a single L2 domain, while current L2 domains are only restricted to several thousands [1]. As a result, current forwarding tables cannot handle such datacenter networks.

The main driver of this expected dramatic growth rate is the deployment of host *virtualization* in datacenter networks. Whereas traditional servers would only use one or two globally-unique MAC addresses per server, contemporary servers use tens of MAC addresses, corresponding to their tens of virtual machines (these generated MAC addresses are not necessarily globally unique anymore). Assuming one MAC address per thread, and considering the joint increase of the numbers of (a) threads per core, (b) cores per processor (socket), (c) sockets per server, and (d) servers per datacenter, it is clear that the product of all these numbers is expected to dramatically increase in the coming years [1]–[3]. In addition

to this expected strong growth in the number of entries per forwarding table, we can also expect a moderate growth in the number of attributes per entry, such as the quality-of-service and security attributes, as well as in the size of these attributes. All in all, we need to be able to handle dramatically larger forwarding tables.

Unfortunately, forwarding tables must be accessed on a per-packet basis. Therefore, they are performance sensitive. Because of the high datacenter rates, they are often implemented in hardware using *on-chip memory*. This makes scaling forwarding tables challenging, because on-chip memory is too small to hold millions of forwarding-table entries in raw form [4], [5].

*The goal of this paper is to study the compressibility of forwarding tables, providing both an analytical framework and experimental insight toward the effective compression of such tables.*

### B. Compressing Forwarding Tables

Forwarding tables, also referred to as Forwarding Information Bases (FIBs), consist of several entries, which associate a *lookup key* to its corresponding *forwarding information*. For instance, the lookup key can be a (VLAN, destination MAC address) pair, and the forwarding information can include the switch target port, the quality-of-service attributes, the security attributes, and additional information used for various modifications that need to be applied to the packet.

Forwarding tables are typically arranged as *two-dimensional tables*. Each row entry in a forwarding table corresponds to a different lookup key, and contains several columns, with the lookup key and forwarding information. Forwarding tables can then be accessed using a hashing mechanism. Upon each incoming packet, the packet lookup key is hashed into one (or several) of the hash table entries. If this entry contains indeed the lookup key, then the forwarding table returns the associated forwarding information.

Forwarding tables are typically redundant, because their fields (columns) can have a highly non-uniform value distribution, and because the values in different fields can be correlated. Therefore, we are interested in studying whether these properties can be used to compress them. However, we would like to keep two key properties that are specific to forwarding tables:

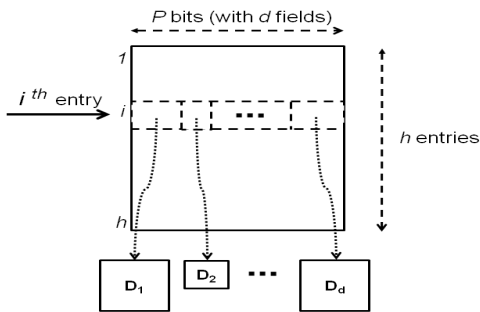


Fig. 1. An illustration of an access to the encoded forwarding table. First, a single encoded entry of (at most)  $P$  bits is retrieved. Then,  $d$  pipelined accesses to  $d$  dictionaries are performed based on the  $d$  binary codewords (of different lengths) composing this encoded entry.

- (i) The ability to access each row entry directly regardless of its row index, and
- (ii) The ability to store each row entry in a fixed-width memory word.

Such criteria prevent or make less attractive traditional encoding techniques. For instance, if the entries are compressed using a Lempel-ZivWelch (LZW) technique and sequentially stored in the memory, it is then difficult to calculate the memory address of a specific entry, and to access it directly based on its entry index.

### C. Overview of the Suggested Solution

As illustrated in Figure 1, we suggest a *novel compression technique for forwarding tables* that supports random accesses.

We propose to encode each of the forwarding-table columns separately, using a dedicated variable-length *binary prefix encoding*, i.e. an encoding in which any codeword is not a prefix of any other codeword. Therefore, each row entry in a  $d$ -column table is represented by the concatenation of  $d$  codewords representing the content of each of its  $d$  fields. Accordingly, its width equals the sum of the  $d$  corresponding codeword lengths. For each row entry we allocate a fixed memory size  $P$ . In addition, each column  $i$  is associated to a dictionary  $D_i$ , which is stored separately in the memory. This dictionary  $D_i$  maps each binary codeword into its represented field value in column  $i$ .

Figure 1 illustrates an access to the encoded forwarding table. At each packet arrival, we first retrieve the  $P$  bits representing the encoding of its  $i^{\text{th}}$  entry. Then, we look up each of the  $d$  dictionaries in a pipelined manner. This way, we sequentially obtain the corresponding values of the  $d$  fields in the  $i^{\text{th}}$  entry after  $d$  pipelined lookups.

Note that the combination of the bounded-width entry property and the prefix property of individual fields yields our two desired properties specified above.

**Example 1.** Table I presents an example of a forwarding table and a possible encoding of that table. There are  $h = 7$  entries in the table, each with  $d = 3$  fields: Target Port, MAC Address and VLAN.

First, Table (A) presents the original table without encoding.

Next, Table (B) illustrates a possible encoding of the forwarding table (the presented spaces in the encoded table are

(A) The original forwarding table

Target Port	MAC Address	VLAN
Te12/1	00:1b:2b:c3:4d:90	Vlan10
Gi11/8	00:00:aa:6c:b1:10	Vlan10
Te12/1	00:00:aa:65:ce:e4	Vlan10
Gi11/24	00:00:aa:65:ce:e4	Vlan200
Gi11/24	00:13:72:a2:a2:0e	Vlan200
Te12/1	00:21:9b:37:7e:14	Vlan10
Gi11/8	00:13:72:a2:a2:0e	Vlan200

(B) Encoded table

0 110 0
10 00 0
0 01 0
11 01 1
11 10 1
0 111 0
10 10 1

(C) The encoding dictionaries

$D_1$ (Target Port)	$D_2$ (MAC Address)	$D_3$ (VLAN)
0 - Te12/1	00 - 00:00:aa:6c:b1:10	0 - Vlan10
10 - Gi11/8	01 - 00:00:aa:65:ce:e4	1 - Vlan200
11 - Gi11/24	10 - 00:13:72:a2:a2:0e	
	110 - 00:1b:2b:c3:4d:90	
	111 - 00:21:9b:37:7e:14	

TABLE I

AN EXAMPLE OF A FORWARDING TABLE WITH  $d = 3$  COLUMNS AND ITS ENCODED REPRESENTATION. (A) SHOWS THE ORIGINAL TABLE WITHOUT ENCODING. THE ENTRY SELECTION MECHANISM IS NOT PRESENTED. (B) ILLUSTRATES THE CORRESPONDING ENCODED FORWARDING TABLE WITH A MAXIMAL ENCODED ENTRY WIDTH OF  $P = 5$  BITS. (SPACES ARE PRESENTED JUST FOR SIMPLICITY AND DO NOT EXIST IN PRACTICE.) (C) DESCRIBES THE  $d = 3$  DICTIONARIES REPRESENTING THE PREFIX ENCODINGS OF THE  $d$  COLUMNS.

just shown to distinguish between the different fields, and do not exist in practice). For instance, the two possible values in the VLAN field are encoded by a simple fixed-length single-bit encoding. In addition, variable-length encodings are used to encode the three distinct values in the Target Port field, as well as the five distinct values of the MAC Address field. With these  $d = 3$  encodings, the obtained maximal width of an encoded entry is  $P = 5$  bits. Thus, in order to support the random-access requirement, we keep  $P = 5$  bits in the memory for each encoded entry, even if it may actually be shorter. For instance, the third entry only needs 4 bits, and therefore wastes 1 bit.

Finally, Table (C) illustrates the encodings used by the  $d = 3$  dictionaries representing the  $d = 3$  columns. The number of entries in each dictionary equals the number of distinct values in the corresponding column in (A).

This example illustrates how prefix encoding enables us to simply concatenate the codewords without a need for any additional separating flags.

In addition, it also shows how the memory size is directly determined by the *worst-case entry width*, i.e. the maximum number of concatenated bits per entry. This is one reason for example why encoding each column using *Huffman coding* [6] would not necessarily address this problem: while it is optimal for the *average* entry width, it may yield a worst case that ends up wasting a significant number of bits. Therefore, for a given forwarding table with  $d$  columns, our goal is to suggest a set of  $d$  prefix encodings that minimizes the maximal encoded entry width. This makes it particularly hard, because each entry width involves all columns, and the maximal one also involves all row entries. Therefore, any strong solution may need to jointly look at all the elements.

#### D. Related Work

The forwarding table scaling problem has been thoroughly analyzed and discussed in networking industry forums and standards organizations, both in the L2 and the L3 layers (e.g. [1], [3]–[5], [7], [8]). These analyses raise the concern that the steep growth rate of forwarding tables has outrun the steady technology growth rate of memory density.

Various creative methods have been proposed to mitigate this scaling problem. First, *FIB aggregation* [9] is an approach that reduces the size of a forwarding table by aggregating table entries with the same values of the next hop field. For instance, two prefix rules that differ only on their last bit can be combined to a single prefix with a shorter prefix. Likewise, a prefix may be eliminated due to redundancy achieved by a shorter prefix with the same next hop. A similar approach is presented in [10]–[13]. For instance, in [11] an optimal representation is found based on dynamic programming where in each rule the value of the next hop can be selected among a short list of possible values. However, this approach assumes that the routing table has a tree structure, i.e. the rules have the form of prefixes in an address field (usually IP address) and cannot be applied in the case of more general fields, such as an exact MAC address with a corresponding attribute.

In addition, there have been studies in other related fields. Clearly, in *information theory*, there are many results on compression and on lower bounds based on entropy [6]. Likewise, *database implementations* also attempt to scale systems by introducing compression techniques [14]–[16]. In particular, in [14] a dictionary-based solution that supports delete and update operations is suggested. A field value is inserted into the dictionary and encoded based on its length and its number of occurrences. Likewise, a block of data is compressed only if its compression ratio is beyond a threshold. However, these studies do not consider the forwarding table restrictions that we mention above, and in particular the bounded entry length. Finally, as mentioned, an orthogonal and complementary consideration is the implementation of the forwarding table using hashing, in order to easily find the row that corresponds to a given key [17]–[19].

#### E. Contributions

*This paper investigates a novel technique for the efficient representation of forwarding tables in datacenter switches, with a support for fast random access and fixed-width memory words.*

We start by studying an offline encoding scheme for a given static FIB table. We prove that it is possible to formulate the encoding problem as an optimization problem, and that the relaxed optimization problem is convex. As a result, we propose a novel algorithm that achieves the optimal worst-case entry width *within a guaranteed constant additive factor of  $d-1$*  in any forwarding table with  $d$  columns, independently of their values and of the number of rows.

Next, we also prove that the general optimization problem is *NP-hard*, as long as there are at least  $d = 7$  columns. We do so by presenting a reduction from the 3 – SAT decision

problem. Therefore, there is little hope for a practical optimal algorithm.

In addition, we also consider the special case of forwarding tables with  $d = 2$  columns, which may only consist of the key and the action. We first suggest tight upper and lower bounds on the optimal encoding width of a two-column table in the general case. We also show a sufficient condition for a fixed-length encoding to be optimal. Moreover, we present an *optimal conditional encoding* of a second column when the encoding of one column is given. Later, we also introduce a bipartite graph that can model the forwarding table, and present additional graph-theoretic bounds on its optimal encoding width. These improved bounds rely on the combinatorial properties of the graph, such as the size of its edge cover.

Finally, we also discuss the limits of our offline model, by analyzing how our compressed forwarding tables may be able to handle *updates*. We show that an additional bit per column and an additional small memory may be sufficient to cope with temporary loads of entry additions, modifications or deletions.

We conclude by evaluating the suggested encoding techniques on real-life forwarding tables as well as on synthetic tables. We demonstrate how our introduced techniques fare against alternative algorithms such as fixed-width and Huffman coding. In particular, we get closer to the compression lower bound, albeit at the cost of an increased complexity.

*Paper Organization:* We start with preliminary definitions in Section II. Then, in Section III we formally define the *FIB encoding problem* and suggest an approximation of its optimal solution. We also show that finding the optimal encoding of a given FIB table is NP-hard in the general case. Next, in Section IV we consider the special case of two-column FIB tables. We present properties of the optimal encoding width and suggest how to find an optimal encoding of a second column given the encoding of the first column. We present additional bounds on the optimal encoding width of a given FIB table in Section V. Later, in Section VI we describe how to support updates of encoded forwarding tables. Last, we evaluate the suggested encoding schemes in Section VII based on synthetic and real-life FIB tables.

## II. MODEL AND PROBLEM FORMULATION

### A. Terminology

We start with the formal definition of terminology used in this paper. For a binary codeword (string)  $x$ , let  $\ell(x)$  denote the length in bits of  $x$ .

**Definition 1** (FIB Table). A FIB table  $A = ((a_1^1, \dots, a_d^1), \dots, (a_1^h, \dots, a_d^h))$  is a two-dimensional matrix of  $h$  rows (entries) and  $d$  columns (fields).

Given a FIB table  $A$ , we denote by  $S_j$  (for  $j \in [1, d]$ ) the set of element values in the  $j^{\text{th}}$  column, i.e.  $S_j = \bigcup_{i=1}^h \{a_j^i\}$ . We also denote by  $n_j$  the number of distinct elements in the  $j^{\text{th}}$  column, i.e.  $n_j = |S_j|$ . Let  $s_{j,1}, \dots, s_{j,n_j}$  denote the elements of  $S_j$ .

For example, in Table I, there are  $n_1 = 3$  values in the first column, and  $S_1 = \{\text{Te12/1}, \text{Gi11/8}, \text{Gi11/24}\}$ .

**Definition 2** (Prefix Encoding). *For a set of elements  $S$ , an encoding  $\sigma$  is an injective mapping  $\sigma : S \rightarrow B$ , where  $B$  is a set of binary codewords of size  $|B| = |S|$ . An encoding is called a prefix encoding if no binary codeword in  $B$  is a prefix (start) of any other binary codeword in  $B$ .*

**Definition 3** (FIB Encoding Scheme). *An encoding scheme  $C_A$  of a FIB Table  $A$  is an ordered set of  $d$  prefix encodings  $C_A = (\sigma_1, \dots, \sigma_d)$ . That is, each  $\sigma_j$  encodes column  $j$ , and is constrained to be a prefix encoding.*

**Definition 4** (FIB Encoding Width). *Let  $C_A = (\sigma_1, \dots, \sigma_d)$  be a FIB encoding scheme of the FIB table  $A$ . The encoding width  $\ell(C_A)$  of  $C_A$  is defined as the maximal sum of the lengths of the  $d$  codewords representing the  $d$  elements in a row of  $A$ , i.e.*

$$\ell(C_A) = \max_{i \in [1, h]} \left( \sum_{j=1}^d \ell(\sigma_j(a_j^i)) \right). \quad (1)$$

**Example 2.** *Let  $A$  be the FIB table from Example 1 with  $h = 7$ ,  $d = 3$ , also described in Table I.(A). Likewise, let  $C_A = (\sigma_1, \sigma_2, \sigma_3)$  be the FIB encoding scheme of  $A$ . Then, as shown in Table I.(C),  $\sigma_1(\text{Te12/1}) = 0$ ,  $\sigma_1(\text{Gi11/8}) = 10$ ,  $\sigma_1(\text{Gi11/24}) = 11$ . As described in Table I.(B),  $\ell(C_A) = \max_{i \in [1, h]} \left( \sum_{j=1}^d \ell(\sigma_j(a_j^i)) \right) = \max(5, 5, 4, 5, 5, 5, 5) = 5$ .*

*Let's compare this scheme with a fixed-length encoding of each column. Since the number of distinct elements in the  $d = 3$  columns of  $A$  are  $n_1 = 3$ ,  $n_2 = 5$  and  $n_3 = 2$ , if  $I_A$  a fixed-length encoding of  $A$  then the elements in the  $j^{\text{th}}$  column are encoded in  $\lceil \log_2 n_j \rceil$  bits, and  $\ell(I_A) = \sum_{j=1}^d \lceil \log_2 n_j \rceil = 2 + 3 + 1 = 6$ .*

### B. Optimal FIB Table Encoding Scheme

For each FIB table  $A$ , we denote by  $OPT(A)$  the optimal encoding width of  $A$ , i.e. the smallest possible encoding width of any encoding scheme of  $A$  such that

$$OPT(A) = \min_{C_A=(\sigma_1, \dots, \sigma_d)} \ell(C_A). \quad (2)$$

Our goal is to find an encoding scheme  $C_A$  that minimizes the encoding width  $\ell(C_A)$ , and therefore reaches this minimum.

## III. GENERAL FIB TABLES

### A. Optimization Problem

We first want to rewrite the problem of finding the optimal encoding scheme for any FIB table  $A$  as an optimization problem.

First, it is well known that the set of codeword lengths in a prefix encoding exists iff it satisfies *Kraft's inequality* [6].

**Property 1** (Kraft's inequality). *There exists a prefix encoding  $\sigma$  of the elements in a set  $S$  with codeword lengths  $\{\ell(\sigma(a)) | a \in S\}$  iff*

$$\sum_{a \in S} 2^{-\ell(\sigma(a))} \leq 1. \quad (3)$$

Therefore, we can now formally present the problem of finding an optimal FIB table encoding scheme for table  $A = ((a_1^1, \dots, a_d^1), \dots, (a_1^h, \dots, a_d^h))$  as the following optimization problem.

$$\begin{aligned} \min \quad & P \\ \text{s.t.} \quad & \sum_{j=1}^d \ell(\sigma_j(a_j^i)) \leq P \quad \forall i \in [1, h] \quad (4a) \end{aligned}$$

$$\sum_{a \in S_j} 2^{-\ell(\sigma_j(a))} \leq 1 \quad \forall j \in [1, d] \quad (4b)$$

$$\ell(\sigma_j(a)) \geq 0 \quad \forall j \in [1, d], \quad \forall a \in S_j \quad (4c)$$

$$\ell(\sigma_j(a)) \in \mathbb{Z} \quad \forall j \in [1, d], \quad \forall a \in S_j \quad (4d)$$

In this optimization problem, we try to minimize the maximal encoding width of a row (denoted here by  $P$ ) while having four sets of constraints. The first set (4a) represents the limitation on the total width of each row. The second set of constraints (4b) requires that each of the encodings  $\sigma_1, \dots, \sigma_d$  should satisfy Kraft's inequality. The last two sets (4c, 4d) illustrate the fact that the codeword length of any element should be a non-negative integer (the constraints (4b) and (4c) together guarantee that if the number of distinct elements in a column is at least two, the codeword lengths are positive). For an optimal solution to this optimization problem, the equality  $P = OPT(A)$  is satisfied.

The optimization problem can be also described as being dependent only on the codeword lengths. To do so, we replace (for  $j \in [1, d], a \in S_j$ ) the value of  $\ell(\sigma_j(a))$  by a corresponding variable. In addition, as explained later, it is easy to derive the sets of codewords given their lengths.

We also denote by the *relaxed FIB encoding problem* the problem achieved by omitting the fourth set of constraints (4d). In this relaxed problem, the codeword lengths are not necessarily integers.

### B. Approximation of the Optimal Encoding

We now show how to find for each FIB table  $A$  of  $d$  columns a solution to the FIB encoding problem that is guaranteed to be within a *fixed additive approximation* of the optimal encoding width  $OPT(A)$ . More specifically, we find an encoding scheme  $C_A$  with encoding width  $\ell(C_A)$  that satisfies  $\ell(C_A) \leq OPT(A) + (d - 1)$ , i.e. its encoding width is larger than the optimal encoding width (in bits) by at most the number of columns in  $A$  minus one. We emphasize that this bound on the additive error of  $(d - 1)$  depends neither on the number of distinct elements  $n_j \leq 2^{W_j}$  (for  $j \in [1, d]$ ) in each of the columns, nor on the number of rows  $h$  in  $A$ .

To obtain such an encoding, we consider the *relaxed FIB encoding problem* defined above, in which the codeword lengths are not necessarily integers. We show that this optimization problem is *convex*, and thus its solution can be found by one of several known algorithms for such problems. We then build as a solution to the original optimization problem, an encoding with codeword lengths achieved by taking the ceiling of each of the codeword lengths in the solution of the relaxed problem.

We show that these new codeword lengths lead to an encoding that satisfies the additive approximation from above.

**Theorem 1.** Let  $(\ell(\bar{\sigma}_j(s_{j,1})), \dots, \ell(\bar{\sigma}_j(s_{j,n_j})))$  (for  $j \in [1, d]$ ) be the codeword lengths of an optimal solution to the relaxed FIB encoding problem. Further, let  $C_A = (\sigma_1, \dots, \sigma_d)$  be an encoding scheme satisfying  $\ell(\sigma_j(s_{j,i})) = \lceil \ell(\bar{\sigma}_j(s_{j,i})) \rceil$  for all  $j \in [1, d], i \in [1, n_j]$ . Then we have:

- (i) The relaxed FIB encoding problem is convex.
- (ii) The encoding width  $\ell(C_A)$  of the encoding scheme  $C_A$  satisfies

$$\ell(C_A) \leq OPT(A) + (d - 1). \quad (5)$$

*Proof:* We first examine the convexity of the relaxed problem. Clearly, its simple objective function is convex. We would like to show that each of the inequality constraint functions (as a function of the codeword lengths) is convex as well. Simply, constraints in the first and the third sets of inequality constraints (4a, 4c) are convex due to their linearity. In addition, we would now like to examine the convexity of the second set of constraints (4b) representing Kraft's inequality. To do so, we define  $d$  functions  $f_1, \dots, f_d$  for the constraints on each of the  $d$  prefix encodings. The function  $f_j$  (for  $j \in [1, d]$ ) receives as a parameter the set of codeword lengths of the elements in  $S_j$ . The function  $f_j(\ell(\sigma_j(s_{j,1})), \dots, \ell(\sigma_j(s_{j,n_j})))$  is defined as  $\sum_{a \in S_j} 2^{-\ell(\sigma_j(a))} = \sum_{i=1}^{n_j} 2^{-\ell(\sigma_j(s_{j,i}))}$ .

We consider two arbitrary encoding schemes  $(\bar{\sigma}_1, \dots, \bar{\sigma}_d), (\hat{\sigma}_1, \dots, \hat{\sigma}_d)$  with codeword lengths  $(\ell(\bar{\sigma}_j(s_{j,1})), \dots, \ell(\bar{\sigma}_j(s_{j,n_j}))), (\ell(\hat{\sigma}_j(s_{j,1})), \dots, \ell(\hat{\sigma}_j(s_{j,n_j})))$  for the elements in  $S_j$  for  $j \in [1, d]$ , respectively. We would like to show that  $f_j(\alpha \cdot \bar{\sigma} + \beta \cdot \hat{\sigma}) \leq \alpha \cdot f_j(\bar{\sigma}) + \beta \cdot f_j(\hat{\sigma})$  for all  $j \in [1, d]$  and  $\alpha, \beta \in [0, 1]$  with  $\alpha + \beta = 1$ . Here, when performing linear operations on the set of codeword lengths, we refer to a new set of codeword lengths, so that each of its codeword lengths is obtained by performing the linear operations on the corresponding length in the original set. Based on the convexity of the function  $g(x) = 2^{-x}$ , we have

$$\begin{aligned} f_j(\alpha \cdot \bar{\sigma} + \beta \cdot \hat{\sigma}) &= \sum_{i=1}^{n_j} 2^{-(\alpha \cdot \ell(\bar{\sigma}_j(s_{j,i})) + \beta \cdot \ell(\hat{\sigma}_j(s_{j,i})))} \\ &\leq \sum_{i=1}^{n_j} \left( \alpha \cdot 2^{-\ell(\bar{\sigma}_j(s_{j,i}))} + \beta \cdot 2^{-\ell(\hat{\sigma}_j(s_{j,i}))} \right) \\ &= \alpha \cdot \sum_{i=1}^{n_j} 2^{-\ell(\bar{\sigma}_j(s_{j,i}))} + \beta \cdot \sum_{i=1}^{n_j} 2^{-\ell(\hat{\sigma}_j(s_{j,i}))} \\ &= \alpha \cdot f_j(\bar{\sigma}) + \beta \cdot f_j(\hat{\sigma}). \end{aligned} \quad (6)$$

Based on these properties, we can finally deduce that the relaxed FIB encoding problem is convex. With this observation, an optimal solution, i.e. a generalized encoding (with non-necessarily integer codeword lengths) can be found by using one of the several known (polynomial time) algorithms for such convex problems (e.g. the ellipsoid algorithm [20]). Let us denote by

$$(\ell(\bar{\sigma}_j(s_{j,1})), \dots, \ell(\bar{\sigma}_j(s_{j,n_j})))$$

the codeword lengths of an optimal solution to the relaxed problem. Let  $OPT_r(A) = \max_{i \in [1, h]} \left( \sum_{j=1}^d \ell(\bar{\sigma}_j(a_j^i)) \right)$  be the maximal encoding width of a row in this optimal solution of the relaxed problem. Again,  $OPT_r(A)$  is not necessarily integer. Clearly, since the set of constraints in the relaxed problem is a subset of the constraints in the FIB encoding problem, the inequality  $OPT_r(A) \leq OPT(A)$  holds. Further, since  $OPT(A)$  is an integer, we have that  $\lceil OPT_r(A) \rceil \leq OPT(A)$ , i.e.  $\lceil OPT_r(A) \rceil$  is a lower bound for  $OPT(A)$ . Clearly, since  $\ell(\sigma_j(s_{j,i})) = \lceil \ell(\bar{\sigma}_j(s_{j,i})) \rceil \geq \ell(\bar{\sigma}_j(s_{j,i}))$ , all the  $d$  constraints representing Kraft's inequality, satisfied by the solution of the relaxed problem, are also satisfied by the set of codeword lengths in  $C_A = (\sigma_1, \dots, \sigma_d)$ .

We now show that the encoding width of the encoding scheme  $C_A$  is within a guaranteed additive approximation of the optimal encoding width of  $A$ ,  $OPT(A)$ . Directly from the definition of  $C_A$ , we can have that

$$\begin{aligned} \ell(C_A) &= \max_{i \in [1, h]} \left( \sum_{j=1}^d \ell(\sigma_j(a_j^i)) \right) = \max_{i \in [1, h]} \left( \sum_{j=1}^d \lceil \ell(\bar{\sigma}_j(a_j^i)) \rceil \right) \\ &< \max_{i \in [1, h]} \left( \sum_{j=1}^d (\ell(\bar{\sigma}_j(a_j^i)) + 1) \right) \\ &= \max_{i \in [1, h]} \left( \sum_{j=1}^d \ell(\bar{\sigma}_j(a_j^i)) \right) + d = OPT_r(A) + d \\ &\leq OPT(A) + d. \end{aligned} \quad (7)$$

Finally, since  $\ell(C_A)$  as well as  $OPT(A)$  are both integers, we can deduce from the (strong) inequality  $\ell(C_A) < OPT(A) + d$  that  $\ell(C_A) \leq OPT(A) + (d - 1)$ . ■

### C. Upper Bound on the Optimal Encoding

We now present a simple upper bound on the optimal encoding of any FIB table with an arbitrary number  $d$  of columns.

Consider a fixed-length encoding scheme  $I_A$  that encodes each column  $j$  using a fixed-length encoding (as seen in Example 2). Then in each column  $j$ , it uses codewords of size  $W_j = \lceil \log_2 n_j \rceil$ , since it needs at least  $n_j > 2^{W_j-1}$  codewords to represent the  $n_j$  different values. As a consequence, the total width of each row in also fixed, and we have

$$\ell(I_A) = \sum_{j=1}^d W_j = \sum_{j=1}^d \lceil \log_2 n_j \rceil. \quad (8)$$

The fixed-length encoding yields an upper bound on the optimal encoding width of the FIB table.

**Property 2.** Let  $A$  be a FIB table of  $d$  columns such that the set of distinct elements in the  $j^{\text{th}}$  column is  $S_j$  with  $|S_j| = n_j \leq 2^{W_j}$ . Then,

$$OPT(A) \leq \sum_{j=1}^d W_j. \quad (9)$$

#### D. NP-HARDNESS

We conclude the section by considering the complexity of finding an optimal encoding scheme of a given FIB table. We show that the above problem is NP-hard for tables with at least seven columns. The proof can be found in the Appendix.

**Theorem 2.** *Given a FIB table  $A$  with  $d \geq 7$  columns, finding an optimal encoding scheme of  $A$  is NP-hard.*

#### IV. TWO-COLUMN FIB TABLES

In datacenters, a popular class of FIB tables are simple L2 MAC tables that contain two columns. The first describes the Target Port, while the second can be viewed as an aggregation of columns representing a collection of other attributes. Indeed, we see in real-life table traces that these additional attributes are hardly ever used. Thus their aggregation has a relatively small set of possible values.

Therefore, from now on, we consider the special case of two-column FIB tables, i.e. FIB tables that satisfy  $d = 2$ . We would like to fundamentally study this case to obtain some intuition on the problem. We show that in this case, in order to find its optimal encoding scheme, a FIB table  $A$  can be simply represented as a bipartite graph. We also suggest an analysis of the optimal encoding width of such two-column FIB tables.

For the sake of simplicity, we assume, unless mentioned otherwise, that the number of distinct elements in each of the  $d = 2$  columns of  $A$  is the same, and that it is a power of two, i.e.  $|S_1| = |S_2| = n = 2^W$ .

##### A. The Optimal Encoding Width of a Two-Column FIB Table

Directly from Property 2, we can deduce that  $OPT(A) \leq 2W$ .

We would like to present an example that shows the tightness of this upper bound by suggesting a two-column FIB table  $A$  that satisfies  $OPT(A) = 2W$ .

Earlier, we present the simple following bound on the codeword lengths in any prefix encoding. Its simple proof follows directly from basic considerations in information theory.

**Lemma 3.** *Let  $\sigma$  be a prefix encoding of the elements of a set  $S$  that satisfies  $|S| = n = 2^W$ . Then, the codeword lengths  $\{\sigma(a) | a \in S\}$  of  $\sigma$  satisfy*

$$\sum_{a \in S} \ell(\sigma(a)) \geq n \cdot W. \quad (10)$$

*In other words, the average length of a codeword of an element is at least  $W$ .*

We now show the tightness of the bound upper bound on  $OPT(A)$ .

**Example 3.** *For  $W \geq 2$  and  $n = 2^W$ , let  $S_1 = \{s_{1,1}, \dots, s_{1,n}\}$  and  $S_2 = \{s_{2,1}, \dots, s_{2,n}\}$ . Consider the two-column FIB table  $A = ((a_1^1, a_2^1), \dots, (a_1^h, a_2^h))$  for  $h = n^2$  such that  $(a_1^{i-n+j}, a_2^{i-n+j}) = (s_{1,i}, s_{2,j})$  for  $i, j \in [1, n]$ . Let  $C_A = (\sigma_1, \sigma_2)$  be a FIB encoding scheme of  $A$ . Clearly,*

*by Lemma 3 ( $\exists i \in [1, n](\ell(\sigma_1(s_{1,i})) \geq W)$ ) and ( $\exists j \in [1, n](\ell(\sigma_2(s_{2,j})) \geq W)$ ). Then,*

$$\begin{aligned} \ell(C_A) &= \max_{k \in [1, h]} \left( \ell(\sigma_1(a_1^k)) + \ell(\sigma_2(a_2^k)) \right) \\ &\geq \ell(\sigma_1(a_1^{i-n+j})) + \ell(\sigma_2(a_2^{i-n+j})) \\ &= \ell(\sigma_1(s_{1,i})) + \ell(\sigma_2(s_{2,j})) \\ &\geq W + W = 2W. \end{aligned} \quad (11)$$

*Since this inequality is correct for any  $C_A$  then  $OPT(A) = \min_{C_A=(\sigma_1, \dots, \sigma_d)} \ell(C_A) \geq 2W$ .*

The next lemma suggests a lower bound on the FIB encoding width of  $A = ((a_1^1, a_2^1), \dots, (a_1^h, a_2^h))$  for  $W \geq 2$ .

**Lemma 4.** *Let  $A$  be a two-column FIB table with  $|S_1| = |S_2| = n = 2^W$  for  $W \geq 2$ . Then, the encoding width of  $A$  satisfies  $W + 2 \leq OPT(A)$ .*

*Proof:* We have to show that any FIB table encoding scheme  $C_A$  of  $A$  satisfies  $\ell(C_A) \geq W + 2$ . Let  $C_A = (\sigma_1, \sigma_2)$  be an arbitrary FIB table encoding scheme of  $A$ . We consider two options for  $C_A$  and show that in both of them  $\ell(C_A) \geq W + 2$ . First, if  $\sigma_1$  is a fixed-length encoding (i.e.  $(\forall a \in S_1)(\ell(\sigma_1(a)) = W)$ ) then consider an element  $b \in S_2$  that satisfies  $\ell(\sigma_1(b)) \geq 2$ . By Kraft's inequality, there must exist such an element when  $W \geq 2$ . Without loss of generality, let  $(a_1^k, a_2^k)$  (for  $k \in [1, h]$ ) be a row such that  $a_2^k = b$ . Then,  $\ell(\sigma_1(a_1^k)) + \ell(\sigma_2(a_2^k)) \geq W + 2$ . Likewise, if  $\sigma_1$  is not a fixed-length encoding, again by Kraft's inequality ( $\exists a \in S_1)(\ell(\sigma_1(a)) \geq W + 1)$ . Let  $(a_1^k, a_2^k)$  (for  $k \in [1, h]$ ) be a row such that  $a_1^k = a$  for such  $a$ . Clearly, for any encoding  $\sigma_2$ ,  $\ell(\sigma_2(a_2^k)) \geq 1$ . Then,  $\ell(\sigma_1(a_1^k)) + \ell(\sigma_2(a_2^k)) \geq W + 2$ . Finally, in both cases

$$\begin{aligned} \ell(C_A) &= \max_{i \in [1, h]} \ell(\sigma_1(a_1^i)) + \ell(\sigma_2(a_2^i)) \\ &\geq \ell(\sigma_1(a_1^k)) + \ell(\sigma_2(a_2^k)) \geq W + 2. \end{aligned} \quad (12)$$

Since the last inequality holds for any encoding scheme  $C_A$ , the result follows. ■

The next example demonstrates that the last lower bound  $OPT(A) \geq W + 2$  is tight for  $W \geq 2$ . To show that, we present, for  $W \geq 2$ , a two-column FIB table  $A$  that satisfies  $OPT(A) = W + 2$ .

**Example 4.** *For  $W \geq 2$  and  $n = 2^W$ , let again  $S_1 = \{s_{1,1}, \dots, s_{1,n}\}$  and  $S_2 = \{s_{2,1}, \dots, s_{2,n}\}$ . Consider the two-column FIB table  $A = ((a_1^1, a_2^1), \dots, (a_1^h, a_2^h))$  for  $h = 2n - 1$  such that  $(a_1^i, a_2^i) = (s_{1,1}, s_{2,i})$  for  $i \in [1, n]$  and  $(a_1^i, a_2^i) = (s_{1,i-(n-1)}, s_{2,n})$  for  $i \in [n, 2n - 1]$ . Clearly,  $S_j = \bigcup_{i=1}^h \{a_j^i\}$  for  $j \in [1, 2]$ . To show that  $OPT(A) = W + 2$ , we suggest a FIB table encoding scheme  $C_A = (\sigma_1, \sigma_2)$  that satisfies  $\ell(C_A) = W + 2$ . Let  $\sigma_1$  satisfy  $\ell(\sigma_1(s_{1,1})) = 1$  and  $\ell(\sigma_1(s_{1,i})) = W + 1$  for  $i \in [2, n]$ . Likewise, let  $\sigma_2$  satisfy  $\ell(\sigma_2(s_{2,i})) = W + 1$  for  $i \in [1, n - 1]$  and  $\ell(\sigma_2(s_{2,n})) = 1$ . Clearly, since  $(n - 1) \cdot 2^{-(W+1)} + 2^{-1} = (2^W - 1) \cdot 2^{-(W+1)} + 2^{-1} < 2^W \cdot 2^{-(W+1)} + 2^{-1} = 1$ , both  $\sigma_1$  and  $\sigma_2$  satisfy Kraft's inequality. Then, for  $i \in [1, n - 1]$  we have that  $\ell(\sigma_1(a_1^i)) + \ell(\sigma_2(a_2^i)) = \ell(\sigma_1(s_{1,1})) + \ell(\sigma_2(s_{2,i})) = 1 + (W + 1) = W + 2$ . We can also see that  $\ell(\sigma_1(a_1^i)) +$*

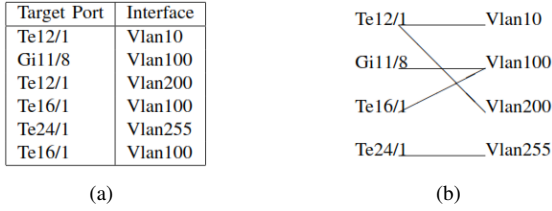


Fig. 2. A two-column forwarding table with  $h = 6$  entries and  $n = 2^W = 4$  (for  $W = 2$ ) distinct elements in each column with (a) the table, and (b) its corresponding bipartite graph. The number of edges in the graph equals the number of distinct entries in the forwarding table.

$\ell(\sigma_2(a_2^n)) = \ell(\sigma_1(s_{1,1})) + \ell(\sigma_2(s_{2,n})) = 1 + 1 = 2$ . Likewise, for  $i \in [n + 1, 2n - 1]$ ,  $\ell(\sigma_1(a_1^i)) + \ell(\sigma_2(a_2^i)) = \ell(\sigma_1(s_{1,i-(n-1)})) + \ell(\sigma_2(s_{2,n})) = (W + 1) + 1 = W + 2$ . We can then have that  $\ell(C_A) = \max_{i \in [1, h]} \left( \sum_{j=1}^2 \ell(\sigma_j(a_j^i)) \right) = W + 2$ . Finally,  $OPT(A) \leq \ell(C_A) = W + 2$  and the lower bound of  $W + 2$  is tight.

Finally, we summarize the lower and upper bounds from Lemma 4 and Property 2 in the following theorem.

**Theorem 5.** *Let  $A$  be a two-column FIB table with  $|S_1| = |S_2| = n = 2^W$  for  $W \geq 2$ . Then, the encoding width of  $A$  satisfies  $W + 2 \leq OPT(A) \leq 2W$ .*

Any two-column FIB table  $A = ((a_1^1, a_2^1), \dots, (a_1^h, a_2^h))$  can be represented by a corresponding bipartite graph  $G_A$  as follows. The two disjoint sets are the sets of distinct elements in each of the two columns. Thus, if an element appears in both of the two columns, it is represented by two different vertices in each of the two disjoint sets. Edges connect elements in the two sets if they appear at least once on the same row of the FIB table. Formally, we define the graph  $G_A = \langle L + R, E \rangle$  such that  $L = S_1$ ,  $R = S_2$  and  $E = \{(x, y) | (\exists i \in [1, h]), (a_1^i, a_2^i) = (x, y)\}$ . Therefore, duplicated rows have no influence on the construction of the bipartite graph, and the graph does not contain parallel edges. It is also easy to see the independence in the order of the rows of the FIB table.

**Example 5.** *Figure 2(a) presents an example of a two-column forwarding table with  $h = 6$  entries and  $n = 2^W = 4$  (for  $W = 2$ ) distinct elements in both columns. The corresponding bipartite graph appears in Figure 2(b). The vertices on the left side of the graph represent the  $n = 4$  distinct elements in the first column, while the vertices on the right side represent the  $n$  distinct elements in the second column. The number of edges in the graph equals the number of distinct entries in the forwarding table.*

Given a FIB encoding scheme  $C_A = (\sigma_1, \dots, \sigma_d)$  of  $A = ((a_1^1, a_2^1), \dots, (a_1^h, a_2^h))$ , we can present its FIB encoding width based on  $G_A$  as  $\ell(C_A) = \max_{(x,y) \in E} \ell(\sigma_1(x)) + \ell(\sigma_2(y))$ . From the construction of  $G_A$ , we can clearly see that the last equation is compatible with Definition 4.

The representation of a FIB table as a bipartite graph can help us to further understand the FIB encoding width based on

tools from graph theory. The next theorem relates the existence of an independent set of a specific size in the bipartite graph  $G_A$  to the value of  $OPT(A)$ .

**Theorem 6.** *Let  $A$  be a two-column FIB table with  $|S_1| = |S_2| = n = 2^W$  and let  $G_A = \langle L + R, E \rangle$  be its corresponding bipartite graph. If there does not exist an independent set of vertices  $U = U_1 \cup U_2$  in  $G_A$ , so that  $U_1 \subseteq L = S_1$ ,  $U_2 \subseteq R = S_2$  and  $|U_1| = |U_2| = \frac{n}{2} + 1$ , then the optimal encoding width of  $A$  necessarily satisfies  $OPT(A) = 2W$ . Namely, the optimal encoding width is achieved by the fixed-length encoding and it cannot be improved by any variable-length encoding.*

*Proof:* Let  $C_A = (\sigma_1, \sigma_2)$  be an arbitrary FIB table encoding scheme of  $A$ . As mentioned earlier, by Property 2,  $OPT(A) \leq 2W$ . We would like to show also that  $\ell(C_A) \geq 2W$  and therefore  $OPT(A) \geq 2W$ . By Kraft's inequality, there are (at least)  $\frac{n}{2} + 1$  elements in  $S_1$  whose codeword lengths in  $\sigma_1$  are at least  $W$  bits. Let  $U_1$  denote the set of their corresponding vertices in  $L$ . Likewise, let  $U_2 \subseteq R$  denote the similar set of vertices representing the elements in  $S_2$  with codeword lengths of at least  $W$  in  $\sigma_2$ . We then have that  $|U_1|, |U_2| \geq \frac{n}{2} + 1$ . If  $\ell(C_A) < 2W$ , then  $(\forall i \in [1, h])$ ,  $\ell(\sigma_1(a_1^i)) + \ell(\sigma_2(a_2^i)) \leq 2W - 1$ . Since  $(\forall x \in U_1, y \in U_2)$ ,  $\ell(\sigma_1(x)) + \ell(\sigma_2(y)) \geq W + W = 2W$ , then necessarily,  $(\forall x \in U_1, y \in U_2)$ ,  $(x, y) \notin E$ . Thus  $U = U_1 \cup U_2$  is an independent set in  $G_A$ . Since  $|U_1|, |U_2| \geq \frac{n}{2} + 1$ , we have a contradiction, and can deduce that indeed  $\ell(C_A) \geq 2W$  and  $OPT(A) = 2W$ . ■

**B. Optimal Conditional Encoding of the Second Column**

We now consider the conditional problem of finding an optimal two-column encoding given that the encoding of the first column is known.

Formally, given a two-column FIB table  $A$  and a known prefix encoding of one of its columns  $\sigma_1 := \bar{\sigma}_1$ , we want to find a prefix encoding  $\sigma_2$  of the second column such that the encoding width  $\ell(C_A)$  of the FIB encoding scheme  $C_A = (\sigma_1 := \bar{\sigma}_1, \sigma_2)$  is minimized. We denote the FIB encoding width of such scheme by  $OPT(A | \sigma_1 := \bar{\sigma}_1)$ , i.e.

$$OPT(A | \sigma_1 := \bar{\sigma}_1) = \min_{C_A = (\sigma_1 := \bar{\sigma}_1, \sigma_2)} \ell(C_A). \quad (13)$$

We would like to suggest an algorithm to find such an optimal conditional encoding. Let  $A$  be a two-column FIB table with two sets of distinct elements in each of the columns  $S_1, S_2$  and an encoding of the first column  $\sigma_1 := \bar{\sigma}_1$ . For  $y \in S_2$ , we denote by  $\phi_{\bar{\sigma}_1}(y)$  the maximal codeword length  $\ell(\bar{\sigma}_1(x))$  of an element  $x \in S_1$  that shares a row with  $y$  in  $A$ .

$$\phi_{\bar{\sigma}_1}(y) = \max\{\ell(\bar{\sigma}_1(x)) | x \in S_1, (x, y) \in E\}. \quad (14)$$

We observe that in order to calculate the optimal encoding  $\sigma_2$  given that  $\sigma_1 := \bar{\sigma}_1$ , it is enough to know  $\phi_{\bar{\sigma}_1}(y)$  for any  $y \in S_2$  while the exact values of  $\ell(\bar{\sigma}_1(s_{1,1})), \dots, \ell(\bar{\sigma}_1(s_{1,n}))$

are not necessarily required. This is because

$$\begin{aligned} OPT(A|\sigma_1 := \bar{\sigma}_1) &= \min_{C_A=(\sigma_1:=\bar{\sigma}_1,\sigma_2)} \ell(C_A) \\ &= \min_{C_A=(\sigma_1:=\bar{\sigma}_1,\sigma_2)} \max_{i \in [1,h]} \left( \ell(\bar{\sigma}_1(a_i^i)) + \ell(\sigma_2(a_i^i)) \right) \\ &= \min_{C_A=(\sigma_1:=\bar{\sigma}_1,\sigma_2)} \max_{y \in S_2} \left( \phi_{\bar{\sigma}_1}(y) + \ell(\sigma_2(y)) \right). \end{aligned} \quad (15)$$

Before presenting the algorithm, we suggest the following lemma.

**Lemma 7.** *There exists a (prefix) encoding  $\bar{\sigma}_2$  of  $S_2$  such that  $(\sigma_1, \sigma_2) = (\bar{\sigma}_1, \bar{\sigma}_2)$  achieves the minimal encoding width of all the FIB encoding schemes  $\{(\sigma_1, \sigma_2) | \sigma_1 := \bar{\sigma}_1\}$  and satisfies*

$$(\forall y \in S_2), \phi_{\bar{\sigma}_1}(y) + \ell(\bar{\sigma}_2(y)) = OPT(A|\sigma_1 := \bar{\sigma}_1). \quad (16)$$

*Proof:* Let  $\hat{C}_A = (\sigma_1 := \bar{\sigma}_1, \hat{\sigma}_2)$  be a FIB encoding scheme that satisfies  $\ell(\hat{C}_A) = OPT(A|\sigma_1 := \bar{\sigma}_1)$ . If  $(\forall y \in S_2), \phi_{\bar{\sigma}_1}(y) + \ell(\hat{\sigma}_2(y)) = OPT(A|\sigma_1 := \bar{\sigma}_1)$  then  $\bar{\sigma}_2 = \hat{\sigma}_2$  can be the requested encoding. Otherwise, by the definition of  $\ell(\hat{C}_A)$ ,  $(\forall y \in S_2), \phi_{\bar{\sigma}_1}(y) + \ell(\hat{\sigma}_2(y)) \leq \ell(\hat{C}_A) = OPT(A|\sigma_1 := \bar{\sigma}_1)$ . For  $y \in S_2$  we use the notation  $q(y) = OPT(A|\sigma_1 := \bar{\sigma}_1) - (\phi_{\bar{\sigma}_1}(y) + \ell(\hat{\sigma}_2(y)))$ . We define an encoding  $\bar{\sigma}_2$  of  $S_2$  based on  $\hat{\sigma}_2$  as  $\bar{\sigma}_2(y) = \hat{\sigma}_2(y) \cdot 0^{q(y)}$  where the  $\cdot$  sign is used to represent the concatenation operation. In particular, if for  $y \in S_2$ ,  $\phi_{\bar{\sigma}_1}(y) + \ell(\hat{\sigma}_2(y)) = OPT(A|\sigma_1 := \bar{\sigma}_1)$  then  $\bar{\sigma}_2(y) = \hat{\sigma}_2(y) \cdot 0^{q(y)} = \hat{\sigma}_2(y)$ . We can see that  $(\forall y \in S_2), \phi_{\bar{\sigma}_1}(y) + \ell(\bar{\sigma}_2(y)) = \phi_{\bar{\sigma}_1}(y) + \ell(\hat{\sigma}_2(y)) + q(y) = OPT(A|\sigma_1 := \bar{\sigma}_1)$ . Thus  $\bar{\sigma}_2$  satisfies the required equality.

To complete the proof we have to show also that  $\bar{\sigma}_2$  is indeed a prefix encoding. We can first verify that the  $n_2 = |S_2|$  codewords of  $\sigma_2$  are different. If for any  $x, y \in S_2$  (s.t.  $x \neq y$ ),  $\bar{\sigma}_2(x) = \hat{\sigma}_2(x) \cdot 0^{q(x)} = \hat{\sigma}_2(y) \cdot 0^{q(y)} = \bar{\sigma}_2(y)$  then necessarily  $\hat{\sigma}_2(x)$  is a prefix of  $\hat{\sigma}_2(y)$  (if  $q(x) \geq q(y)$ ) or  $\hat{\sigma}_2(y)$  is a prefix of  $\hat{\sigma}_2(x)$  (if  $q(x) \leq q(y)$ ), a contradiction since  $\hat{\sigma}_2$  is a prefix encoding.  $\bar{\sigma}_2$  also preserves the prefix property since if  $\bar{\sigma}_2(x) = \hat{\sigma}_2(x) \cdot 0^{q(x)}$  is a prefix of  $\bar{\sigma}_2(y) = \hat{\sigma}_2(y) \cdot 0^{q(y)}$  then again either  $\hat{\sigma}_2(x)$  is a prefix of  $\hat{\sigma}_2(y)$  or  $\hat{\sigma}_2(y)$  is a prefix of  $\hat{\sigma}_2(x)$ . ■

Based on the last lemma we can deduce from Kraft's inequality, the exact value of  $OPT(A|\sigma_1 := \bar{\sigma}_1)$ .

**Theorem 8.** *The optimal conditional FIB encoding width of  $A$  satisfies*

$$OPT(A|\sigma_1 := \bar{\sigma}_1) = \left\lceil \log_2 \left( \sum_{y \in S_2} 2^{\phi_{\bar{\sigma}_1}(y)} \right) \right\rceil. \quad (17)$$

*Proof:* We first would like to show that  $OPT(A|\sigma_1 := \bar{\sigma}_1) \geq \left\lceil \log_2 \left( \sum_{y \in S_2} 2^{\phi_{\bar{\sigma}_1}(y)} \right) \right\rceil$ . Consider a FIB encoding  $(\sigma_1, \sigma_2) = (\bar{\sigma}_1, \bar{\sigma}_2)$  that has an encoding width of  $OPT(A|\sigma_1 := \bar{\sigma}_1)$  and  $(\forall y \in S_2)$  satisfies  $\phi_{\bar{\sigma}_1}(y) + \ell(\bar{\sigma}_2(y)) = OPT(A|\sigma_1 := \bar{\sigma}_1)$ . The last equality can be presented also as  $\ell(\bar{\sigma}_2(y)) = OPT(A|\sigma_1 := \bar{\sigma}_1) - \phi_{\bar{\sigma}_1}(y)$ . Such an encoding exists by Lemma 7. By Kraft's inequality

$\bar{\sigma}_2$  satisfies

$$\begin{aligned} 1 &\geq \sum_{y \in S_2} 2^{-\ell(\bar{\sigma}_2(y))} = \sum_{y \in S_2} 2^{-\left( OPT(A|\sigma_1 := \bar{\sigma}_1) - \phi_{\bar{\sigma}_1}(y) \right)} \\ &2^{OPT(A|\sigma_1 := \bar{\sigma}_1)} \geq \sum_{y \in S_2} 2^{\phi_{\bar{\sigma}_1}(y)} \\ OPT(A|\sigma_1 := \bar{\sigma}_1) &\geq \left\lceil \log_2 \left( \sum_{y \in S_2} 2^{\phi_{\bar{\sigma}_1}(y)} \right) \right\rceil, \end{aligned} \quad (18)$$

where the last inequality is obtained since  $OPT(A|\sigma_1 := \bar{\sigma}_1)$  is of course an integer. To complete the proof we would also like to show that  $OPT(A|\sigma_1 := \bar{\sigma}_1) \leq \left\lceil \log_2 \left( \sum_{y \in S_2} 2^{\phi_{\bar{\sigma}_1}(y)} \right) \right\rceil$  by finding a FIB encoding scheme  $(\sigma_1, \sigma_2) = (\bar{\sigma}_1, \bar{\sigma}_2)$  whose encoding width satisfies this upper bound. To do so, we denote by  $v_{\bar{\sigma}_1}$  the term  $\left\lceil \log_2 \left( \sum_{y \in S_2} 2^{\phi_{\bar{\sigma}_1}(y)} \right) \right\rceil$ . Then,

$$\begin{aligned} \sum_{y \in S_2} 2^{-(v_{\bar{\sigma}_1} - \phi_{\bar{\sigma}_1}(y))} &= \sum_{y \in S_2} 2^{\phi_{\bar{\sigma}_1}(y) - v_{\bar{\sigma}_1}} = \sum_{y \in S_2} 2^{\phi_{\bar{\sigma}_1}(y)} \cdot 2^{-v_{\bar{\sigma}_1}} \\ &\leq \sum_{y \in S_2} 2^{\phi_{\bar{\sigma}_1}(y)} \cdot \left( \sum_{y \in S_2} 2^{\phi_{\bar{\sigma}_1}(y)} \right)^{-1} = 1. \end{aligned} \quad (19)$$

Thus there exists a prefix encoding  $\bar{\sigma}_2$  with codeword lengths  $\ell(\bar{\sigma}_2(y)) = (v_{\bar{\sigma}_1} - \phi_{\bar{\sigma}_1}(y))$  for  $y \in S_2$ . The FIB encoding scheme  $C_A = (\bar{\sigma}_1, \bar{\sigma}_2)$  satisfies

$$\begin{aligned} \ell(C_A) &= \max_{y \in S_2} \left( \phi_{\bar{\sigma}_1}(y) + \ell(\bar{\sigma}_2(y)) \right) \\ &= \max_{y \in S_2} \left( \phi_{\bar{\sigma}_1}(y) + (v_{\bar{\sigma}_1} - \phi_{\bar{\sigma}_1}(y)) \right) \\ &= \max_{y \in S_2} \left( v_{\bar{\sigma}_1} \right) = v_{\bar{\sigma}_1}. \end{aligned} \quad (20)$$

Then,  $OPT(A|\sigma_1 := \bar{\sigma}_1) \leq v_{\bar{\sigma}_1} = \left\lceil \log_2 \left( \sum_{y \in S_2} 2^{\phi_{\bar{\sigma}_1}(y)} \right) \right\rceil$ . ■

Now, when the codeword lengths of the elements of  $S_2$  are known from Lemma 7 and Theorem 8, it is easy to find such an encoding. To do so, we can find the codewords of the elements of  $S_2$  in an increasing order of their required lengths. For each element  $y \in S_2$ , we allocate an arbitrary codeword  $\bar{\sigma}_2(y)$  in the required length such that none of the previously allocated codewords is a prefix of the current codeword. It is easy to verify that Kraft's inequality guarantees that such allocation is always possible.

## V. LOWER BOUNDS ON THE OPTIMAL ENCODING WIDTH OF TWO-COLUMN FIB TABLES

We would like now to suggest additional lower bounds on the optimal encoding width of a given two-column FIB table  $A$  as defined above. We can calculate this bound without solving the relaxed FIB encoding problem.

Let again  $G_A = \langle L + R, E \rangle = \langle S_1 + S_2, E \rangle$  with  $|S_1| = |S_2| = n = 2^W$ . We first recall an additional definition from graph theory.

**Definition 5 (Edge Cover).** *For an undirected graph  $G = \langle V, E \rangle$ , a set of edges  $S \subseteq E$  is called an edge cover of  $G$  if*



$\bigcup_{(x,y) \in S} \{x, y\} = V$ , i.e. each of the vertices in  $V$  is incident on at least one of the edges in  $S$ .

Clearly,  $G_A$  does not include any isolated vertices, i.e. any vertex is connected to at least one vertex and the value it represents appears in at least one of the rows of  $A$ . Therefore, an edge cover always exists. For a graph  $G$ , we denote by  $\rho(G)$  the edge covering number of  $G$ , i.e. the minimal number of edges in an edge cover of  $G$ . We first show a simple property of  $\rho(G_A)$ .

**Lemma 9.** *The edge covering number  $\rho(G_A)$  of  $G_A$  satisfies  $n \leq \rho(G_A) \leq 2 \cdot (n - 1)$ .*

*Proof:* Consider the edges in a minimal edge cover in an arbitrary order. The first edge covers two vertices. Any additional edge adds one or two vertices to the set of vertices covered by previous edges in the edge cover. Since there are  $2n$  vertices in  $G_A$  then clearly  $n \leq \rho(G_A) \leq 2n - 1$ . To show that  $\rho(G_A) \leq 2n - 2$ , we show that there must be two independent edges in  $G_A$ , i.e. edges that do not share any vertex. Consider an edge  $(s_{1,i}, s_{2,j})$  for  $i, j \in [1, n]$ . If all the other edges in  $E$  share at least one vertex with this edge, we must have that  $E = \{(s_{1,i}, s_{2,k}) | k \in [1, n]\} \cup \{(s_{1,k}, s_{2,j}) | k \in [1, n]\}$ . Then, for instance the edges  $(s_{1,i}, s_{2,k}), (s_{1,k}, s_{2,j})$  for  $k \neq i, j$  are independent. Finally, the two independent edges can contribute two covered vertices to the edge cover and necessarily  $\rho(G_A) \leq 2n - 2 = 2 \cdot (n - 1)$ . ■

Let  $\alpha = \rho(G_A)/n$ . The next theorem suggests a lower bound on the optimal encoding width of  $A$  based on the value of  $\rho(G_A)$ .

**Theorem 10.** *The optimal encoding width of  $A$  satisfies*

$$OPT(A) \geq \left\lceil \frac{2 \cdot (W + \alpha - 1)}{\alpha} \right\rceil. \quad (21)$$

*Proof:* For  $G_A = \langle S_1 + S_2, E \rangle$ , let  $S \subseteq E$  be an edge cover of  $G_A$  of size  $\rho(G_A) = \alpha \cdot n$ . Let  $C_A = (\sigma_1, \sigma_2)$  be an arbitrary FIB table encoding scheme of  $A$ . We will show that

$$\max_{(x,y) \in S} \ell(\sigma_1(x)) + \ell(\sigma_2(y)) \geq \left\lceil \frac{2 \cdot (W + \alpha - 1)}{\alpha} \right\rceil. \quad (22)$$

Thus, since  $S \subseteq E$  then necessarily

$$\ell(C_A) = \max_{(x,y) \in E} \ell(\sigma_1(x)) + \ell(\sigma_2(y)) \geq \left\lceil \frac{2 \cdot (W + \alpha - 1)}{\alpha} \right\rceil. \quad (23)$$

To do so, we consider the subset of the rows of  $A$  represented by the edge cover  $S$ . Let  $D = ((d_1^1, d_2^1), \dots, (d_1^{\rho(G_A)}, d_2^{\rho(G_A)}))$  represent these rows, s.t.  $|D| = \rho(G_A) = \alpha \cdot n$ . By Definition 5 of the edge cover,  $S_1 = \bigcup_{i=1}^{\rho(G_A)} \{d_1^i\}$  and  $S_2 = \bigcup_{i=1}^{\rho(G_A)} \{d_2^i\}$ , i.e. all the elements of  $S_1$  appear at least once in the first column of the rows of  $D$  as well as all the elements of  $S_2$  in the second. Let  $X_D(C_A)$  denote the total sum of codeword lengths of the elements in  $D$  using the encoding  $C_A = (\sigma_1, \sigma_2)$ , s.t.  $X_D(C_A) = \sum_{i=1}^{\rho(G_A)} (\ell(\sigma_1(d_1^i)) + \ell(\sigma_2(d_2^i)))$ .

By manipulating Kraft's inequality, we can have that the total sum of codeword lengths of the  $n$  distinct elements in each column is at least  $n \cdot W$ . Of course, the other  $\alpha \cdot n - n =$

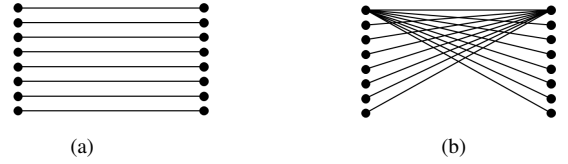


Fig. 3. Two examples of bipartite graphs with  $n = 2^W = 8$  (for  $W = 3$ ) vertices in each graph side representing two two-column FIB tables  $A_1$  and  $A_2$ . In the first graph (a), there is a perfect match and the size of the (minimal) edge cover is  $8 = n = \alpha \cdot n$  with  $\alpha = 1$  and  $OPT(A_1) = 2W = 6$ . In the second (b), the size of the edge cover is  $14 = \alpha \cdot n$  with  $\alpha = 1.75$  and  $OPT(A_2) \geq 5$ .

$(\alpha - 1) \cdot n$  elements in each column are encoded in at least one bit each. Thus,

$$\begin{aligned} X_D(C_A) &= \sum_{i=1}^{\rho(G_A)} (\ell(\sigma_1(d_1^i)) + \ell(\sigma_2(d_2^i))) \\ &= \sum_{i=1}^{\rho(G_A)} \ell(\sigma_1(d_1^i)) + \sum_{i=1}^{\rho(G_A)} \ell(\sigma_2(d_2^i)) \\ &\geq n \cdot W + (\alpha - 1) \cdot n + n \cdot W + (\alpha - 1) \cdot n \\ &= 2n \cdot W + 2(\alpha - 1) \cdot n. \end{aligned} \quad (24)$$

We recall that the total number of rows is  $\alpha \cdot n$  and thus there is at least one row in  $D$  that is encoded in at least  $\left\lceil \frac{X_D(C_A)}{\alpha \cdot n} \right\rceil$  bits. Thus the encoding width of  $C_A$  satisfies

$$\begin{aligned} \ell(C_A) &= \max_{(x,y) \in E} \ell(\sigma_1(x)) + \ell(\sigma_2(y)) \\ &\geq \max_{(x,y) \in S} \ell(\sigma_1(x)) + \ell(\sigma_2(y)) \geq \left\lceil \frac{X_D(C_A)}{\alpha \cdot n} \right\rceil \\ &\geq \left\lceil \frac{2n \cdot W + 2(\alpha - 1) \cdot n}{\alpha \cdot n} \right\rceil \\ &= \left\lceil \frac{2 \cdot (W + \alpha - 1)}{\alpha} \right\rceil \geq \left\lceil \frac{2W}{\alpha} \right\rceil. \end{aligned} \quad (25)$$

Finally, Since  $C_A$  is an arbitrary encoding scheme, we can deduce the result. ■

**Example 6.** Figure 3 illustrates the two corresponding bipartite graphs for two-column FIB tables  $A_1$  and  $A_2$ , with  $n = 2^W = 8$  (for  $W = 3$ ) distinct elements in each column.

For  $A_1$ , the bipartite graph includes a perfect match and the size of the edge cover is  $8 = n = \alpha \cdot n$  with  $\alpha = 1$ . Theorem 10 implies that  $OPT(A_1) \geq \left\lceil \frac{2 \cdot (W + \alpha - 1)}{\alpha} \right\rceil = 2W$ .

For  $A_2$ , the size of the minimal edge cover is  $14 = 2 \cdot (n - 1) = \alpha \cdot n$  with  $\alpha = 1.75$ . Thus  $OPT(A_2) \geq \left\lceil \frac{2 \cdot (3 + 1.75 - 1)}{1.75} \right\rceil = 5$ . Further, by encoding the first element in a column in a single bit and others in four bits, we have that indeed  $OPT(A_2) = 5$ .

The above lower bound can be slightly improved in some cases. We first recall a well known property from graph theory. In a bipartite graph of two pairs of  $n$  vertices, there exists a minimal edge cover of size  $\alpha \cdot n$  iff there exists a maximum matching of size  $2n - \alpha \cdot n = (2 - \alpha) \cdot n$ . Let  $\beta = (2 - \alpha)$ . To improve the last presented lower bound, we calculate the average encoding of a the  $\beta \cdot n$  rows in a maximum matching. Let  $k = \lfloor \log_2(\beta \cdot n) \rfloor$  s.t.  $2^k \leq \beta \cdot n < 2^{k+1}$ .

**Theorem 11.** *The optimal encoding width of  $A$  satisfies*

$$OPT(A) \geq \left\lceil \frac{2 \cdot \left( (k+2) \cdot (\beta \cdot n) - 2^{k+1} + I(1-\beta) \right)}{\beta \cdot n} \right\rceil \quad (26)$$

where  $I$  is defined as  $I(1-\beta) = 1$  when  $(1-\beta) > 0$  and  $I(1-\beta) = 0$  otherwise.

*Proof:* Let  $B_1, B_2$  denote the elements of the first and second column in the rows of such a maximum matching, respectively. Then,  $B_1 \subseteq S_1, B_2 \subseteq S_2$  and  $|B_1| = |B_2| = \beta \cdot n$ . Consider again an arbitrary encoding scheme  $C_A = (\sigma_1, \sigma_2)$ .

We first assume that  $I(1-\beta) = 1$  (and accordingly  $\beta \cdot n < n$ ). The lengths of the codewords of the elements in  $B_1, B_2$  satisfy

$$\sum_{x \in B_1} 2^{-\ell(\sigma_1(x))} < 1, \quad \sum_{y \in B_2} 2^{-\ell(\sigma_2(y))} < 1. \quad (27)$$

(Here, an equality cannot hold since other elements in  $S_1 \setminus B_1, S_2 \setminus B_2$  also have positive codeword lengths.) When satisfying these strong inequalities, the minimal total sum of the codeword lengths of the elements in  $B_1$  (as well as in  $B_2$ ) is achieved when, among the  $\beta \cdot n \in [2^k, 2^{k+1})$  elements, exactly  $2^{k+1} - \beta \cdot n - 1$  elements have a codeword length of  $k$  while the other  $\beta \cdot n - (2^{k+1} - \beta \cdot n - 1) = 2 \cdot \beta \cdot n - 2^{k+1} + 1$  elements have a codeword lengths of  $k+1$ . Thus, the sum of codeword length of the elements in  $B_1$  (as well as in  $B_2$ ) is at least

$$\begin{aligned} & k \cdot \left( 2^{k+1} - \beta \cdot n - 1 \right) + (k+1) \cdot \left( 2 \cdot \beta \cdot n - 2^{k+1} + 1 \right) \\ &= (k+2) \cdot (\beta \cdot n) - 2^{k+1} + 1. \end{aligned} \quad (28)$$

Since each row of the  $\beta \cdot n$  rows in the maximum matching includes a single element from  $B_1$  and a single element from  $B_2$ , the average encoding width of these rows is at least

$$\frac{2 \cdot \left( (k+2) \cdot (\beta \cdot n) - 2^{k+1} + 1 \right)}{\beta \cdot n} \quad (29)$$

and the result follows.

Likewise, if  $I(1-\beta) = 0$  then  $\beta \cdot n = n = 2^W$  and  $k = W$ . The minimal sum of codeword lengths is achieved when all the elements has a codeword length of  $W$ . Thus the total encoding length of the elements in  $B_1$  (as well as in  $B_2$ ) is at least

$$\begin{aligned} W \cdot \beta \cdot n &= k \cdot 2^k = (k+2) \cdot 2^k - 2^{k+1} \\ &= (k+2) \cdot (\beta \cdot n) - 2^{k+1}. \end{aligned} \quad (30)$$

■

**Example 7.** *Let  $A$  be a FIB table satisfying  $W = 4$  and  $n = 2^W = 16$ , with edge covering number  $\rho(G_A) = 20 = \alpha \cdot n$  for  $\alpha = 1.25$ . Theorem 10 implies that  $OPT(A) \geq \left\lceil \frac{2 \cdot (W + \alpha - 1)}{\alpha} \right\rceil = \left\lceil \frac{2 \cdot (4 + 1.25 - 1)}{1.25} \right\rceil = \lceil 6.8 \rceil = 7$ .*

*In this case, the size of a maximum matching is  $\beta \cdot n = (2 - \alpha) \cdot n = 0.75 \cdot n = 12$ . Here,  $\beta = 0.75$  and  $I(1-\beta) = 1$  where  $2^k \leq \beta \cdot n < 2^{k+1}$  for  $k = 3$ . In each column, an encoding with a minimal total encoding length*

*of the  $\beta \cdot n = 12$  elements might have encoding lengths of 3, 3, 3, 4, 4, 4, 4, 4, 4, 4, 4, 4 with a total encoding length of*

$$\begin{aligned} 3 \cdot 3 + 9 \cdot 4 &= 45 = (3+2) \cdot 12 - 2^{3+1} + 1 \\ &= (k+2) \cdot (\beta \cdot n) - 2^{k+1} + I(1-\beta). \end{aligned} \quad (31)$$

*Then, by Theorem 11 we have an improved lower bound of*

$$OPT(A) \geq \left\lceil \frac{2 \cdot 45}{\beta \cdot n} \right\rceil = \left\lceil \frac{2 \cdot 45}{12} \right\rceil = \lceil 7.5 \rceil = 8. \quad (32)$$

## VI. SUPPORTING UPDATES

*Forwarding tables need to support updates.* These updates can include an insertion of a new entry, a deletion of an existing entry, or a modification of some of the fields of an entry. For instance, they may follow the addition or deletion of VMs in the datacenter, as well VM migrations. In this section, we discuss how such updates can be supported.

Let  $P$  be the fixed allocated number of bits for the encoding of each entry. To support such updates, we assume that each entry includes an additional bit that indicates its validity.

In general, the update process includes two steps: *an immediate step required for coherency, and another procedure that can occur offline to improve performance.* This second step can be run, for instance, either periodically in time or after a fixed number of updates. Meanwhile, in addition to the encoded table, we make use of a small dedicated memory for keeping a small number of entries.

Dealing with a change in an entry or with an insertion can be done in a similar way. We describe several possible scenarios. First, *if all the elements in the updated entry appear in the current dictionaries*, we simply encode a new entry based on these dictionaries. If its total width is not greater than  $P$ , no further changes are required. If it is greater than  $P$ , we set the original entry (in case of a change) as invalid and temporarily add the new entry (without encoding it) to the small memory.

*If the updated entry includes a new element*, we try to associate it with a codeword. In a specific column, we can do so only if the sum that appears in Kraft's inequality is smaller than one. Further, the minimal possible codeword length depends on the value of this sum. For the sake of simplicity, we suggest to allocate for each new element a codeword with the minimal possible length. If we cannot allocate any new codeword in one of these columns (due to an equality in the corresponding Kraft's inequality), or if the allocated codewords yield an entry that is too long, the updated entry is again added to the small memory.

Dealing with *an entry deletion* is easy and requires setting the entry as invalid. Such a change does not require any changes in the encodings of other entries, since the maximal width of the rest of the encoded entries is clearly bounded by the maximal width before change. Reduction in the maximal width of the table after an entry removal may be achieved by running the compression algorithm as part of the offline procedure.

The offline process includes the calculation of an efficient encoding for the existing encoded entries and for the unencoded entries stored in the small memory. Later, this memory is cleared. The value of  $P$  may be changed, as well as the

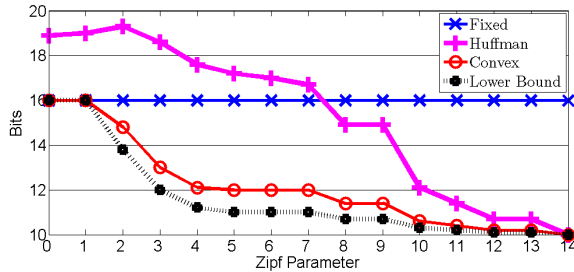
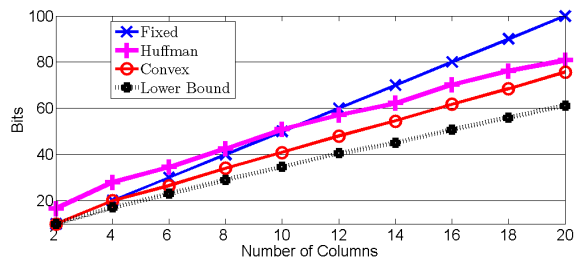
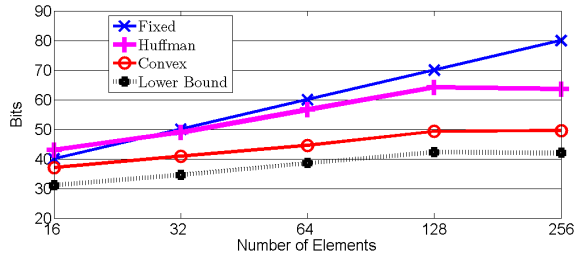


Fig. 4. Effectiveness of the encoding algorithm on synthetic FIB tables with  $d = 2$  columns and  $W = 8$ , as a function of the Zipf parameter. Elements are selected according to the Zipf distribution with the given parameter. By Theorem 5, the optimal encoding width is at least  $W + 2 = 10$  and at most  $2W = 16$ . Since  $d = 2$ , our suggested *convex* scheme takes at most  $d - 1 = 1$  bit more than the theoretical lower bound. In comparison, the fixed-length encoding has a fixed width of  $2W = 16$  bits, while the Huffman encoding sometimes achieves a greater width.



(a) Effectiveness of the encoding algorithm as a function of the number of columns with  $2^W = 32$  distinct elements per column.



(b) Effectiveness of the encoding algorithm as a function of the number of distinct elements per column with  $d = 10$  columns.

Fig. 5. Analysis of the encoding algorithm as a function of table parameters.

encoded version of existing entries. It is possible to endow the table encoding algorithm with better update capabilities by choosing an encoding with a slightly larger value of  $P$  that guarantees short codeword lengths for possible future inserted elements. Meanwhile, between two offline processes, any lookup should consider the encoded entries in the FIB and the unencoded entries in the small dedicated memory.

## VII. EXPERIMENTAL RESULTS

We now turn to conducting experiments to examine the distribution of the optimal encoding width and the performance of our suggested encoding algorithm from Section III-B. In the experiments, we rely on both synthetic and real-life FIB tables.

### A. Effectiveness on Synthetic FIB tables

We first evaluate our suggested scheme on *synthetic two-column FIB tables*. Each table has 1000 entries with  $2^W = 256$  (for  $W = 8$ ) distinct elements in each of the  $d = 2$  columns. To generate the table, we first make sure that each element appears at least once in each column by adding an entry containing it with an additional random element in the other column. This additional random element is selected according to a Zipf distribution of parameter  $s$ . Then, we add additional entries with  $d = 2$  random elements selected according to the same distribution. Intuitively, a Zipf distribution with a lower parameter  $s$  is closer to the uniform distribution.

First, in Fig. 4, we compare our suggested scheme, denoted “convex”, with the fixed-length and Huffman encoding schemes. The results are based on the average of 10 simulations. We present the lower bound given by the fractional solution of the relaxed optimization problem (Section III-B). We also show our suggested scheme, which takes the ceiling of the fractional solution of the relaxed problem, unless it exceeds  $2W$ , in which case it simply uses fixed-length encoding. As suggested by Theorem 1, our scheme is always within  $d - 1 = 1$  bit of the fractional-optimum lower bound. For instance, for a Zipf parameter  $s = 4$ , the widths for the Huffman, fixed-length, and our suggested encodings are respectively 17.6,  $2W = 16$ , and 12.1, while the lower bound is 11.2. More generally, the lower the Zipf parameter, the closer the lower bound to  $2W$ .

Next, in Fig. 5, we plot the encoding efficiency as a function of the number of columns (a) and of the number of elements appearing in each column (b). In both plots we take  $h = 1000$  entries and a Zipf parameter  $s = 2$ . In Fig. 5(a) the number of elements is 32, i.e.  $W = 5$ , and in Fig. 5(b) the number of columns is  $d = 10$ . We can see that in (a), Huffman becomes better for the worst-case entry width. Intuitively, it is because the ratio of the standard deviation of the sum of the encodings of  $d$  columns to the expected sum decreases with  $d$  as  $\frac{1}{\sqrt{d}}$ .

We want to stress that while our scheme performs better, it is also much slower to run than the simpler schemes. As mentioned in the implementation section, it is destined to offline or background use. For instance, in Fig. 5(a), an unoptimized implementation of our scheme on a basic Intel i5 core respectively took 0.82 seconds and 23.1 seconds for the two-column and ten-column tables, while the Huffman scheme ran in the respectively- negligible times of 7 and 11 milliseconds.

### B. Real-Life FIB tables

We also conduct experiments on three typical real-life enterprise-network tables. The tables are collected from three different switches of three different switch vendors, each switch using a different number of columns. All tables were collected in enterprise networks, the first in the United States with 2790 entries and the other in China with 903 and 428 entries. For each table, we present the size of the original raw table without compression. We compare it to the total size of the compressed table, including the dictionary for three compression algorithms: the fixed-length encoding, the

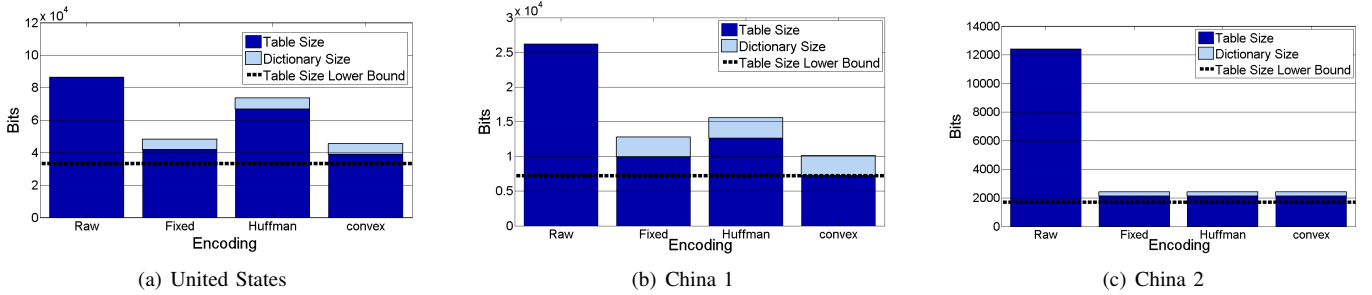


Fig. 6. Total memory size of encoded real-life tables including the dictionary size. All forwarding tables were collected in enterprise networks, in the United States as well as in China.

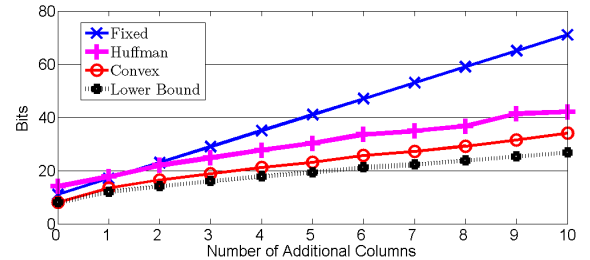
Huffman encoding and our suggested encoding. The results are presented in Fig. 6.

For instance, in the second real-life table (China 1), presented in Fig. 6 (b), we consider the 903 entries of two columns representing the VLAN and the Target-Port fields. Without encoding, each entry is allocated 29 bits per entry and thus the raw data without encoding requires  $903 \cdot 29 = 26187$  bits. Using the three encodings, the size of the dictionaries is almost fixed and equals approximately 2900 bits. An entry width of 11 and 14 bits is achieved in the fixed-length encoding and in the Huffman encoding, respectively, while in the proposed encoding we achieve an entry width of only 8 bits. This leads to an improvement of 20.7%, 34.9%, and 61.3% in the total memory size (including the dictionary) compared to the fixed-length, Huffman and Raw solutions, respectively.

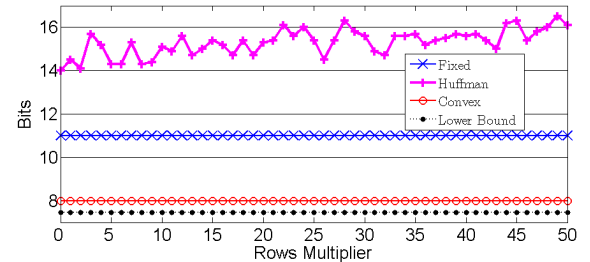
### C. Extrapolated Tables

We want to extrapolate the behaviors of these schemes from real-life limited data sets to the larger expected data sets to be encountered in future datacenter network architectures. The approach we take herein is to amplify the available data sets by natural operations that maintain the general properties of the data. We carry out these amplifications to scale both the number of columns and the number of rows in the table. For each of the two, we demonstrate the performance of our proposed algorithm in comparison to both the fixed-width and Huffman encodings. The plots are the result of averaging 10 simulations.

**Additional Columns:** We extend the China 1 table by duplicating and appending the port column several times. In each new column, we permute the appearances of the elements uniformly at random. The motivation is to model additional forwarding attributes that would follow a distribution that is similar to an existing one. Rearranging the row locations of the duplicate elements adds randomness to the extended file. As seen in Fig. 7(a), for each additional column the fixed encoding increases by a constant, as expected since the number of unique elements in each column does not vary. The Huffman and convex encodings increase at a slower pace, thereby increasing the encoding effectiveness as the number of columns grows. Of course, since this table is extended synthetically, we urge caution in over-interpreting the results.



(a) Effectiveness of the encoding algorithm on the extended table, as a function of the number of columns.



(b) Effectiveness of the encoding algorithm on the extended table, as a function of the number of rows.

Fig. 7. Effectiveness of the encoding algorithm on the extended table (based on the China 1 table).

**Additional Rows:** We also extend the China 1 table by uniformly choosing a row from the table, duplicating it and inserting it at the end of the table. This process is repeated until the desired number of rows is reached. The motivation for this operation is to model additional forwarding entries that share common instances with existing ones. As in the previous case, the random choice of duplicated rows adds realistic randomness to the extended file. As seen in Fig. 7(b), the fixed encoding remains the same as no new unique elements were added. Likewise, the convex encoding also remains unaltered due to no new constraints being added to the convex problem. However, the Huffman encoding's length exhibits high variation in addition to slightly increasing with the number of rows. The Huffman encoding differs from the other two because it depends on the appearance frequencies of the elements.

In summary, the convex encoding becomes more efficient as the number of columns grows, and remains stable as the number of rows increases in the above scenario. Of course,

since these table are extended synthetically, we urge caution in over-interpreting the results.

### VIII. CONCLUSION

In this paper, we saw how datacenter virtualization is causing a dramatic rise in the number of entries in forwarding tables, such that it becomes impossible to implement forwarding tables in current memory sizes without any compression algorithm. We then investigated the compressibility of forwarding tables, by introducing a novel forwarding table architecture with separate encoding in each column. Our architecture supports fast random accesses and fixed-width memory words. We also later explained how our architecture can handle table updates.

Later, we suggested an encoding whose per-entry memory requirement is guaranteed to be within a small additive constant of the optimum. Finally, we evaluated this new scheme against the fixed-width and Huffman schemes, using both synthetic and real-life forwarding tables from different switches, switch vendors, and network country locations. As future work, we plan to check how using additional hardware such as CAM and TCAM can help in further compressing the forwarding tables.

### IX. ACKNOWLEDGMENT

We would like to thank Yishay Mansour and Haim Kaplan for their helpful participation and suggestions.

This work was partly supported by the European Research Council Starting Grant No. 210389, by a European Commission Marie Curie CIG grant, by the Israel Science Foundation grant No. 1241/12, by the German-Israeli Foundation for Scientific Research and Development, by the Intel ICRI-CI Center, by the Hasso Plattner Center for Scalable Computing and by the Israel Ministry of Science and Technology. Ori Rottenstreich is the Google Europe Fellow in Computer Networking and a Jacobs-Qualcomm Fellow.

### REFERENCES

- [1] I. Gashinsky, "Datacenter scalability panel," in *North American Network Operators Group, NANOG 52*, 2011.
- [2] R. N. Mysore *et al.*, "Portland: a scalable fault-tolerant layer 2 data center network fabric," in *SIGCOMM*, 2009.
- [3] T. Narten *et al.*, "Problem statement for ARMD," IETF, draft-ietf-armd-problem-statement, work in progress, 2012.
- [4] G. Hankins, "Pushing the limits, a perspective on router architecture challenges," in *North American Network Operators Group, NANOG 53*, 2011.
- [5] D. Meyer, L. Zhang, and K. Fall, "Report from the IAB workshop on routing and addressing," IETF, RFC 4984, 2007.
- [6] T. M. Cover and J. A. Thomas, *Elements of information theory (2. ed.)*. Wiley, 2006.
- [7] T. Narten, "On the scalability of internet routing," IETF, draft-narten-radir-problem-statement, work in progress, 2010.
- [8] P. Francis, X. Xu, H. Ballani, D. Jen, R. Raszuk, and L. Zhang, "FIB suppression with virtual aggregation," IETF, draft-ietf-grow-va, work in progress, 2011.
- [9] X. Zhao, Y. Liu, L. Wang, and B. Zhang, "On the aggregatability of router forwarding tables," in *INFOCOM*, 2010.
- [10] R. Draves, C. King, V. Srinivasan, and B. Zill, "Constructing optimal IP routing tables," in *INFOCOM*, 1999.
- [11] Q. Li, D. Wang, M. Xu, and J. Yang, "On the scalability of router forwarding tables: Nexthop-selectable FIB aggregation," in *Infocom Mini-Conference*, 2011.

- [12] Y. Liu, X. Zhao, K. Nam, L. Wang, and B. Zhang, "Incremental forwarding table aggregation," in *GLOBECOM*, 2010.
- [13] Z. A. Uzmi, M. Nebel, A. Tariq, S. Jawad, R. Chen, A. Shaikh, J. Wang, and P. Francis, "Smalta: practical and near-optimal FIB aggregation," in *CoNEXT*, 2011.
- [14] M. Pöss and D. Potapov, "Data compression in Oracle," in *VLDB*, 2003.
- [15] R. Johnson, V. Raman, R. Sidle, and G. Swart, "Row-wise parallel predicate evaluation," *PVLDB*, vol. 1, no. 1, pp. 622–634, 2008.
- [16] K. Stolze, V. Raman, R. Sidle, and O. Draese, "Bringing BLINK closer to the full power of SQL," in *BTW*, 2009.
- [17] Y. Kanizo, D. Hay, and I. Keslassy, "Optimal fast hashing," in *IEEE Infocom*, 2009.
- [18] A. Kirsch and M. Mitzenmacher, "The power of one move: Hashing schemes for hardware," *IEEE/ACM Trans. Netw.*, 2010.
- [19] Y. Kanizo, D. Hay, and I. Keslassy, "Hash tables with finite buckets are less resistant to deletions," *Computer Networks*, 2012.
- [20] Y. Nesterov and A. Nemirovskii, *Interior-point polynomial algorithms in convex programming*. Society for Industrial Mathematics, 1987, vol. 13.

### APPENDIX

To prove Theorem 2, we will classically show a *reduction from the 3 – SAT decision problem* with a correct number of variables to the optimal encoding problem of a table with 7 columns.

In the 3 – SAT decision problem, we would like to determine whether a CNF (conjunctive normal form) formula with three variables per clause is satisfiable. Let  $k$  be any number, and suppose we are given a 3 – SAT formula with  $n$  variables, where  $n = \frac{2^k - 1}{3}$ . Denote by  $x_1, \dots, x_n$  the variables of the formula. We will show that the formula is satisfiable *iff* some carefully constructed FIB table with seven columns has encoding width of at most  $3k + 6$ . Since  $k = \log(3n + 1)$ , having  $2^{O(k)}$  symbols in the table means that the problem size is polynomial in  $n$ .

In addition, to differentiate between the symbols in different columns, we will use the letters  $A, B, C, D, E, F, G$  for the columns. For example,  $A_3$  will be a symbol in the first column. The subscript 3 means that in our reduction  $A_3$  represents the variable  $x_3$  in the first column. In addition, we also define a *dummy blank* symbol in each column, with index 0 (e.g.,  $A_0$  in the first column), as well as some similar *special* symbol of index  $-1$ .

We begin by presenting the table  $M$  that has to be encoded. We start with an empty  $M$ , and successively add several entries. The construction is a bit involved, and we explain it later.

- 1) For every  $1 \leq j \leq 2^{3k-1}$ , add the two entries

$$(A_0, B_0, C_0, D_0, E_0, F_j, G_0) \quad (33a)$$

$$(A_0, B_0, C_0, D_0, E_j, F_0, G_0). \quad (33b)$$

- 2) For every  $1 \leq i \leq n$  and every  $1 \leq j \leq 2^{2k-1}$ , add the following six entries:

$$(A_i, B_0, C_0, D_0, E_0, F_0, G_j) \quad (34a)$$

$$(\bar{A}_i, B_0, C_0, D_0, E_0, F_0, G_j) \quad (34b)$$

$$(A_0, B_i, C_0, D_0, E_0, F_0, G_j) \quad (34c)$$

$$(A_0, \bar{B}_i, C_0, D_0, E_0, F_0, G_j) \quad (34d)$$

$$(A_0, B_0, C_i, D_0, E_0, F_0, G_j) \quad (34e)$$

$$(A_0, B_0, \bar{C}_i, D_0, E_0, F_0, G_j). \quad (34f)$$

- 3) For every  $1 \leq i \leq n$  and every  $1 \leq j \leq 2^k$ , add the following six entries:

$$(A_i, \bar{B}_i, C_0, D_j, E_0, F_0, G_0) \quad (35a)$$

$$(\bar{A}_i, B_i, C_0, D_j, E_0, F_0, G_0) \quad (35b)$$

$$(A_i, B_0, \bar{C}_i, D_j, E_0, F_0, G_0) \quad (35c)$$

$$(\bar{A}_i, B_0, C_i, D_j, E_0, F_0, G_0) \quad (35d)$$

$$(A_0, B_i, \bar{C}_i, D_j, E_0, F_0, G_0) \quad (35e)$$

$$(A_0, \bar{B}_i, C_i, D_j, E_0, F_0, G_0). \quad (35f)$$

- 4) For every  $1 \leq j \leq 2^{2k-1}$ , add the following special three entries:

$$(A_{-1}, B_0, C_0, D_0, E_0, F_0, G_j) \quad (36a)$$

$$(A_0, B_{-1}, C_0, D_0, E_0, F_0, G_j) \quad (36b)$$

$$(A_0, B_0, C_{-1}, D_0, E_0, F_0, G_j). \quad (36c)$$

- 5) Finally, for every clause  $C$  that contains the literals  $x_i, x_j$  and (say)  $\bar{x}_k$  we will have a single entry of the form

$$(A_i, B_j, \bar{C}_k, D_0, E_0, F_0, G_0) \quad (37)$$

This also applies equally to clauses with 0, 2 or 3 negated variables. The order of the variables will not matter for us.

Let's now provide some intuition for the naming and ideas behind the reduction, before proving it formally. Note that only entries of the form of Equation (37) actually depend on the input. The rest of the entries are used to reduce from 3-SAT to encoding matrices.

The symbols  $S_0$ , with  $S \in \{A, B, C, D, E, F, G\}$ , represent dummy blanks, and must be encoded with just one character in an encoding scheme of width  $3k + 6$ .

The symbols  $S_i$  and  $\bar{S}_i$  where  $S \in \{A, B, C\}$  represent the literals  $x_i$  and  $\bar{x}_i$  respectively. We will later explain why in an encoding scheme of width  $3k + 6$  (if one exists) it must be the case that one of  $S_i, \bar{S}_i$  is encoded by a codeword of length  $k$ , and the other has length  $k + 1$ . We will interpret the one which has codeword length  $k$  as the value of the truth assignment. In entries of the form of Equation (37), we essentially require that at least one of the literals of the clause has a short codeword, which would correspond to the literal being satisfied.

The following lemma proves the easy direction of the result:

**Lemma 12.** *If the formula is satisfiable, there is an encoding scheme with an encoding width of  $3k + 6$ .*

*Proof:* Let  $T$  be some truth assignment. We present the encoding scheme  $C_M = (\sigma_1, \dots, \sigma_d)$  (for  $d = 7$ ) that satisfies  $\ell(C_M) = 3k + 6$ .

- 1) Each blank symbol  $S_0$  will be encoded by one bit, for any  $S \in \{A, B, \dots, G\}$ .
- 2) Each symbol  $E_j, F_j$  will be encoded by  $3k$  bits. Note that there are  $2^{3k-1}$  such symbols.
- 3) Each symbol  $G_j$  will be encoded by  $2k$  bits. Note that there are  $2^{2k-1}$  such symbols.
- 4) Each symbol  $D_j$  will be encoded with  $k + 1$  bits. Note that there are  $2^k$  such symbols.

- 5) Each special symbol  $A_{-1}, B_{-1}, C_{-1}$  will be encoded by  $k + 1$  bits.

- 6) For every  $S \in \{A, B, C\}$ , if  $T(x_i)$  is true then  $S_i$  is encoded with  $k$  bits and  $\bar{S}_i$  is encoded in  $k + 1$  bits. In addition, if  $T(x_i)$  is false then  $S_i$  is encoded with  $k + 1$  bits and  $\bar{S}_i$  is encoded in  $k$  bits.

We can verify that each of the  $d = 7$  encodings of  $C_M$  satisfy Kraft's inequality. For instance, in the first three columns we have (for  $j \in [1, 3]$  and respectively  $S \in \{A, B, C\}$ )  $2^{-\ell(\sigma_j(S_{-1}))} + 2^{-\ell(\sigma_j(S_0))} + \sum_{i=1}^n \left( 2^{-\ell(\sigma_j(S_i))} + 2^{-\ell(\sigma_j(\bar{S}_i))} \right) = 2^{-(k+1)} + 2^{-1} + n \cdot (2^{-k} + 2^{-(k+1)}) = 2^{-(k+1)} + 2^{-1} + (2^k - 1) \cdot 2^{-(k+1)} = 1$ .

In addition, the encodings width of all entries is at most  $3k + 6$ . For example, since in a truth assignment at least one literal is satisfied, we have that each entry in the last set of entries (37) is encoded in at most  $2 \cdot (k + 1) + k + 4 = 3k + 6$  bits. ■

We would like to suggest another property showing the other direction. To do so, we present several lemmas.

**Lemma 13.** *Suppose that there is an encoding scheme of width  $3k + 6$ . Then every symbol  $S_0$  has to be encoded with just one bit, for every  $S \in \{A, B, \dots, G\}$ .*

*Proof:* There are  $2^{3k-1}$  symbols of the form  $E_j$  and  $F_j$  (in addition to  $E_0, F_0$ ). By Kraft's inequality, at least one of these symbols (in both columns) is encoded by (at least)  $3k$  bits. Encoding them in entries (33a) and (33b) requires that all  $S_0$  symbols (including  $E_0$  and  $F_0$ ) have a codeword of 1 bit. ■

**Lemma 14.** *Each symbol  $S_{-1}, S_i$  and  $\bar{S}_i$  is encoded by at most  $k + 1$  bits, for every  $S \in \{A, B, C\}$  and  $1 \leq i \leq n$ .*

*Proof:* There are  $2^{2k-1}$  symbols of the form  $G_j$  (in addition to  $G_0$ ). At least one of them is encoded in (at least)  $2k$  bits. Consider entries (34a)-(34f) and (36a)-(36c). Even if all the  $S_0$  are encoded with one bit (as we know they are), if any  $S_i$  is encoded with more than  $k + 1$  bits the obtained encoding width will be larger than  $3k + 6$ . ■

Now, again by Kraft's inequality, we must have that (for  $S \in \{A, B, C\}$ ), besides  $S_0$  that is encoded in a single bit, at most  $n$  symbols out of  $\{S_i, \bar{S}_i, S_{-1}\}$  have codeword lengths of at most  $k$ , and at least  $n + 1$  signals have codeword lengths of  $k + 1$ . This is because otherwise,  $2^{-1} + (n + 1) \cdot 2^{-k} + n \cdot 2^{-(k+1)} = 2^{-1} + (3n + 2) \cdot 2^{-(k+1)} = 2^{-1} + (2^k + 1) \cdot 2^{-(k+1)} > 1$ .

We now add some definitions. For  $S \in \{A, B, C\}$ , define:  $L_S = \{i : \text{both } S_i \text{ and } \bar{S}_i \text{ have codeword lengths of at most } k\}$ ,  $H_S = \{i : \text{both } S_i \text{ and } \bar{S}_i \text{ have codeword lengths of } k + 1\}$ , and  $T_S = \{i : \text{exactly one of } S_i \text{ and } \bar{S}_i \text{ has codeword length of at most } k\}$ .

**Lemma 15.** *In any encoding scheme of width  $\leq 3k + 6$  we have  $H_S = L_S = \emptyset$  for every  $S \in \{A, B, C\}$ .*

*Proof:* By the upper bound from above on the number of symbols encoded in at most  $k$  bits, we have  $|H_S| \geq |L_S|$ . For every  $i \in H_A$ , both  $A_i$  and  $\bar{A}_i$  have codeword lengths

of  $k + 1$ . But in this case entries (35a) - (35f) require that  $B_i, \bar{B}_i, C_i$  and  $\bar{C}_i$  all have codeword lengths of at most  $k$ , that is  $i \in L_B, L_C$ . By a similar argument for  $H_B$  and  $H_C$  we get that  $2|H_A| + 2|H_B| + 2|H_C| = L_A + L_B + L_C$  which can only be true if  $|H_S| = |L_S| = 0$  for every  $S \in \{A, B, C\}$ . ■

The following lemma lets us derive a truth assignment.

**Lemma 16.**  *$A_i$  has a codeword length of  $k$  iff  $B_i$  has a codeword length of  $k$  iff  $C_i$  has a codeword length of  $k$ .*

*Proof:* By contradiction. Suppose that  $A_i$  has a codeword length of  $k$  but  $B_i$  has a codeword length of  $k + 1$ . By Lemma 15,  $\bar{A}_i$  also has a codeword length of  $k + 1$ . There are  $2^k$  symbols of the form  $D_j$  (in addition to  $D_0$ ), thus at least one of them is encoded in at least  $k + 1$  bits. Then, the entry of the form (35b) with this symbol and  $\bar{A}_i, B_i$  must have a width of at least  $3 \cdot (k + 1) + 4 = 3k + 7$ . Contradiction to the encoding width of the scheme. The rest of the cases follow similarly. ■

We can now elicit a truth assignment which satisfies the formula:  $T(i) = true$  if and only if  $A_i$  has an a codeword length  $k$ .

**Lemma 17.** *The assignment  $T$  satisfies the formula.*

*Proof:* Consider some clause. As the corresponding entry (37) has an encoding width of at most  $3k + 6$ , at least one of three symbols in this entry must be encoded in  $k$  bits. The corresponding literal will satisfy the clause under the truth assignment  $T$ . ■

Directly from the last lemma, we can deduce the requested property.

**Property 3.** *Suppose that there is an encoding scheme of width  $3k + 6$ . Then the original formula is satisfiable.*

Finally, we can present the proof of Theorem 2 by combining both lemmas.

*Proof of Theorem 2:* Combining Lemma 12, Lemma 3, and the hardness of the 3-SAT decision problem, we deduce the hardness of finding the optimal encoding width of a FIB table with  $d = 7$ , and therefore the hardness of finding an optimal encoding scheme as well. To show the hardness for a table  $M'$  with  $d > 7$  columns, we can simply add to a table  $M$  with 7 columns,  $d - 7$  additional columns, each with exactly two distinct elements. Since all the elements in the additional columns can be encoded in a single bit we must have that  $OPT(M') = (3k + 6) + (d - 7)$  if and only if  $OPT(M) = 3k + 6$ . ■