

A tight upper bound on the size of the antidictionary of a binary string

Hiroyoshi Morita^{1†} and Takahiro Ota²

¹Graduate School of Information Systems, University of Electro-Communications, Chofugaoka 1-5-1, Chofu, Tokyo, 182-8585 Japan. ²Dept. of Electronic Engineering, Nagano Prefectural Institute of Technology, Shimonogo 813-8, Ueda, Nagano, 386-1211 Japan.

A tight upper bound of the size of the antidictionary of a binary string is presented. And it is shown that the size of the antidictionary of a binary string is always smaller than or equal to that of its dictionary. Moreover, an algorithm to reconstruct its dictionary from its antidictionary is given.

Keywords: antidictionary, minimum forbidden words, suffix trees, data compression, ECG

1 Introduction

An antidictionary is a set of words that never appear in a binary string. In 2000, Crochemore et al. (2000) presented a compression algorithm of binary text using antidictionary called DCA. Their coding algorithm has been tested on the Calgary Corpus, and their experimental results show that we get compression ratios equivalent to those of most common compressors such as pkzip. Recently, an online source coding scheme based on DCA is presented to apply for compressing losslessly ECG (ElectroCardioGram) in Ota and Morita (2004). Experimental results show that their algorithm achieved 10% smaller compression ratio than LZ ones.

In this article, we present an upper bound of the size of the antidictionary of a binary string. The upper bound we obtained is stronger than that in Crochemore et al. (1998). And it is tight in the sense there exists a string to attain the bound. We also proved that the antidictionary of a binary string is always smaller than or equal to that of the dictionary of the same string. Moreover, we give an algorithm to reconstruct the dictionary from the antidictionary.

This article is organized as follows. Section 2 gives definitions on antidictionary with some examples. In Sections 3 and 4, we investigate the size of the antidictionary of a given string and derive a tight upper bound on its size. Section 5 presents an algorithm to reconstruct the dictionary from the antidictionary of a given string and Section 6 summarizes our results.

2 Definitions on Antidictionary

Let \mathcal{A} be the binary alphabet $\{0, 1\}$ and \mathcal{A}^* be the set of all finite-length strings over \mathcal{A} including the null string of length zero, denoted by λ . The dictionary $\mathcal{D}(\mathbf{x})$ of a binary string $\mathbf{x} = x_1x_2\dots x_n \in \mathcal{A}^*$ is defined as the set of all the substrings of \mathbf{x} :

$$\mathcal{D}(\mathbf{x}) = \{x_\ell x_{\ell+1} \dots x_m \mid 1 \leq \ell \leq m \leq n\} \cup \{\lambda\}.$$

For example, if $\mathbf{x} = 01011$, then $\mathcal{D}(\mathbf{x})$ is given by

$$\mathcal{D}(01011) = \{\lambda, 0, 1, 01, 10, 11, 010, 011, 101, 0101, 1011, 01011\}.$$

Letting $\mathcal{D}^c(\mathbf{x}) = \mathcal{A}^* \setminus \mathcal{D}(\mathbf{x})$, a string $\mathbf{v} = v_1v_2\dots v_m \in \mathcal{D}^c(\mathbf{x})$ such that

$$v_1v_2\dots v_{m-1} \in \mathcal{D}(\mathbf{x}) \text{ and } v_2v_3\dots v_m \in \mathcal{D}(\mathbf{x})$$

Tab. 1: LIST OF ANTIDictionary OF SEVERAL x 'S.

x	$\mathcal{AD}(x)$	x	$\mathcal{AD}(x)$
0	{00}	000	{0000}
1	{11}	001	{000, 10, 11}
00	{000}	010	{00, 101, 11}
01	{00, 10, 11}	011	{00, 111, 10}

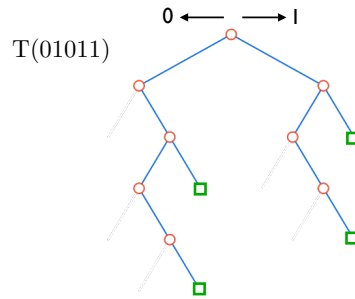


Fig. 1: Suffix trie of $x = 01011$.

is called a minimal forbidden word (MFW) of x . The antidictionary of x , denoted by $\mathcal{AD}(x)$, is defined as the set of all the MFW's of x . In case of $x = 01011$, $\mathcal{AD}(x) = \{00, 110, 111, 1010\}$. Table 1 shows the antidictionaries of several binary strings.

Let $\mathcal{S}(x)$ be the set of suffices of x :

$$\mathcal{S}(x) = \{\lambda\} \cup \{x_i x_{i+1} \dots x_n \mid 1 \leq i \leq n\}.$$

The suffix trie $T(x)$ is a tree structure (Gusfield, 1997) such that every suffix of x is stored as a path from the root to a node in $T(x)$ where every edge is labeled with a symbol in \mathcal{A} . Figure 1 shows the suffix trie of $x = 01011$. Note that some suffices are implicitly represented as paths from the root to internal nodes in $T(x)$. In fact, every string in $\mathcal{D}(x)$ can be represented as a path from the root to a node. Reversely, for every node p in $T(x)$, a string represented by a path from the root to p is in $\mathcal{D}(x)$. Hence, we obtain the following statement.

Proposition 1 (Suffix Trie representing Dictionary). *A node in $T(x)$ corresponds uniquely to an element in $\mathcal{D}(x)$ and vice versa.*

3 A Necessary Condition on MFW's

A necessary condition that a string in \mathcal{A}^* is an MFW of x can be derived by adding new nodes to $T(x)$ as follows: If p is a leaf in $T(x)$, then create two nodes connected to p . Otherwise, and if p has only a child node, create a new node connected to p . The obtained tree, denoted by $T_{\text{ex}}(x)$, is a binary tree such that every internal node has two child nodes (See Fig. 2). Moreover, let $w(r)$ be the string associated with the path from the root to a node r in $T_{\text{ex}}(x)$.

It is shown that every MFW of x is represented by a path from the root to a leaf p in $T_{\text{ex}}(x)$ since for its parent node q , $w(q) \in \mathcal{D}(x)$ but $w(p) \notin \mathcal{D}(x)$ and there exists $a \in \mathcal{A}$ such that $w(p) = w(q)a$. Hence, we obtain the following proposition.

Proposition 2. *If $v \in \mathcal{A}^*$ is an MFW of x , then there exists a leaf p in $T_{\text{ex}}(x)$ such that $v = w(p)$.*

From Propositions 1 and 2, we can derive a rough estimation on the size of $\mathcal{AD}(x)$. Throughout this article, the size, or the cardinality of a set \mathcal{S} is denoted by $\#\mathcal{S}$.

Theorem 1. *For $x \in \mathcal{A}^*$, we have*

$$\#\mathcal{AD}(x) \leq \#\mathcal{D}(x) + 1.$$

[†]This work was supported by Grant-in-Aid for Scientific Research no.173601782332 of Japan Society for the Promotion of Science.

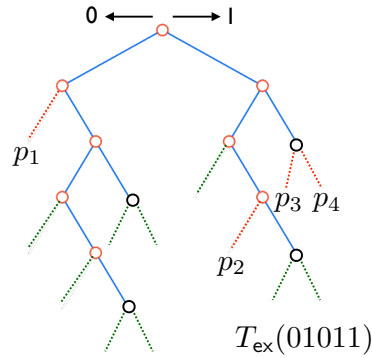


Fig. 2: $T_{\text{ex}}(\mathbf{x})$ corresponding to the suffix trie in Figure 1.

Proof. Let m_k be the number of nodes having k child nodes in $T(\mathbf{x})$ where $0 \leq k \leq 2$. From Proposition 1, we have

$$\#\mathcal{D}(\mathbf{x}) = m_0 + m_1 + m_2.$$

Since $m_0 = m_2 + 1$, we can rewrite it as

$$\#\mathcal{D}(\mathbf{x}) = 2m_0 + m_1 - 1.$$

On the other hand, the number of leaves in $T_{\text{ex}}(\mathbf{x})$ is upper bounded by $2m_0 + m_1$. Hence Proposition 2 gives

$$\#\mathcal{AD}(\mathbf{x}) \leq 2m_0 + m_1.$$

Therefore, we have $\#\mathcal{AD}(\mathbf{x}) \leq \#\mathcal{D}(\mathbf{x}) + 1$. □

Recently, Janson et al. (2004) investigated the average size of a dictionary of a random binary string generated by a mixing model. It is asymptotically equal to $n^2/2$. In the same paper, they did more precise analysis on the average behavior of the number of distinct substrings in a string of length n over an alphabet of size A . The size of an antidictionary, however, is much smaller than that of the dictionary on average as we will discuss below.

4 A Tight Upper Bound on Size of Antidictionary

Hereafter, for any node r in $T_{\text{ex}}(\mathbf{x})$, let $\sigma(r)$ be a node such that $w(\sigma(r))$ is equal to the string obtained by removing the first symbol of $w(r)$. The following theorem gives a necessary and sufficient condition that $w(p)$ for a leaf p in $T_{\text{ex}}(\mathbf{x})$ is an MFW of \mathbf{x} .

Theorem 2 (A necessary and sufficient condition on MFW). *For a leaf node p in $T_{\text{ex}}(\mathbf{x})$, $w(p)$ is an MFW of \mathbf{x} if and only if $\sigma(p)$ is an internal node in $T_{\text{ex}}(\mathbf{x})$.*

Proof. Suppose that p is a leaf in $T_{\text{ex}}(\mathbf{x})$ and q is its parent node. Then there exists $a \in \mathcal{A}$ such that $w(p) = w(q)a$. From the definition of $T_{\text{ex}}(\mathbf{x})$, $w(p) \notin \mathcal{D}(\mathbf{x})$ while $w(q) \in \mathcal{D}(\mathbf{x})$. If $\sigma(p)$ is an internal node in $T_{\text{ex}}(\mathbf{x})$, then $w(\sigma(p)) \in \mathcal{D}(\mathbf{x})$ from Proposition 1. Moreover, there exists $b \in \mathcal{A}$ such that $w(p) = bw(\sigma(p))$. Hence $w(p)$ is an MFW of \mathbf{x} .

Conversely, assume that $w(p)$ is an MFW of \mathbf{x} for a leaf p in $T_{\text{ex}}(\mathbf{x})$. Rewriting $w(p)$ as cu with a certain symbol $c \in \mathcal{A}$, string u corresponds to node $\sigma(p)$, that is, $u = w(\sigma(p))$. Since $w(p)$ is an MFW, $w(\sigma(p)) \in \mathcal{D}(\mathbf{x})$. Therefore, $\sigma(p)$ is an internal node in $T_{\text{ex}}(\mathbf{x})$. □

Corollary 1. *Suppose that p is a leaf in $T_{\text{ex}}(\mathbf{x})$ and its parent node q is an internal node in $T(\mathbf{x})$. Then, $w(p)$ is an MFW of \mathbf{x} if and only if node $\sigma(q)$ has two child nodes in $T(\mathbf{x})$.*

Proof. If node $\sigma(q)$ has two child nodes, one of them is $\sigma(p)$. Hence $\sigma(p)$ is a node in $T(\mathbf{x})$. Thus, it is also an internal node in $T_{\text{ex}}(\mathbf{x})$. Conversely, assume that $\sigma(p)$ is an internal node in $T_{\text{ex}}(\mathbf{x})$. Since p is a leaf in $T_{\text{ex}}(\mathbf{x})$, its parent node q has two child nodes including p . Hence, $\sigma(q)$ does so too. □

Corollary 1 shows that the size of $\mathcal{AD}(\mathbf{x})$ is at least m_2 where m_2 is the number of nodes having two child nodes in $T(\mathbf{x})$ as defined in Theorem 1. In Figure 2, strings 00, 1010, 110 and 111 are corresponding to leaves p_1 to p_4 , respectively. Since their parents satisfy the conditions of Corollary 1, these strings are MFW's of $\mathbf{x} = 01011$. And there are no other leaves whose parents do so.

Theorem 3 (MFW sprouting from leaves in $T(\mathbf{x})$). *For a leaf q in $T(\mathbf{x})$, string $w(p)$ associated with p that is one of q 's child nodes in $T_{\text{ex}}(\mathbf{x})$ is an MFW if and only if the path from the root to q is the shortest one among all the leaves in $T(\mathbf{x})$.*

Proof. First, we assume that q is the leaf q^* with the shortest path in $T(\mathbf{x})$. Then $\sigma(q^*)$ is an internal node in $T(\mathbf{x})$. Thus, $\sigma(q^*)$ has at least one child node r in $T(\mathbf{x})$. Since there exists a child node p^* of q^* such that $r = \sigma(p^*)$, $w(p^*)$ is an MFW of \mathbf{x} .

Conversely, if $q \neq q^*$, then $w(q^*)$ is a suffix of $w(q)$ since $w(q)$ is strictly longer than $w(q^*)$ and both of them are suffixes of \mathbf{x} . Let g be a child node of q in $T_{\text{ex}}(\mathbf{x})$ such that $w(p^*)$ is a suffix of $w(g)$ where $w(p^*)$ is an MFW defined above. Thus neither $w(g)$ is in $\mathcal{D}(\mathbf{x})$. Therefore any suffixes of $w(g)$ that are longer than or equal to $w(p^*)$ are not in $\mathcal{D}(\mathbf{x})$. Hence, $w(g)$ is not an MFW. Taking the contraposition completes the proof of Theorem 3. \square

For example, two leaves associated with strings 110 and 111 in $T_{\text{ex}}(01011)$ of Figure 2 are MFW's that satisfy the condition of Theorem 3. Finally, we have that $\mathcal{AD}(01011) = \{00, 110, 1010, 111\}$.

Theorem 4 (An improved bound of Theorem 1). *Given a binary string \mathbf{x} of length n , we have*

$$\#\mathcal{AD}(\mathbf{x}) \leq n + 1.$$

And if \mathbf{x} is the all-one string $1 \cdots 1$ or the all-zero string $0 \cdots 0$, then

$$\#\mathcal{AD}(\mathbf{x}) = 1.$$

Proof. Combining the results of Corollary 1 and Theorem 3, we have

$$\#\mathcal{AD}(\mathbf{x}) \leq m_2 + 2.$$

Since $m_0 \leq n$ and $m_0 = m_2 + 1$, the above inequality is evaluated further from above as follows:

$$\#\mathcal{AD}(\mathbf{x}) \leq m_0 + 1 \leq n + 1.$$

If \mathbf{x} is the all-one string $1 \cdots 1$ of length n , the all-one string of length $n + 1$ is an MFW of \mathbf{x} and any other strings are not. Therefore, $\#\mathcal{AD}(\mathbf{x}) = 1$. In case \mathbf{x} is the all-zero string of length n , the same argument derives the equality. \square

Since the equality holds for $\mathbf{x} = 01$ (see Table 1), the upper bound obtained in Theorem 4 is tight. In case of the binary alphabet, Corollary 9 in Crochemore et al. (1998) is translated into

$$\#\mathcal{AD}(\mathbf{x}) \leq \begin{cases} 3 & \text{if } |\mathbf{x}| \leq 2, \\ 2 & \text{else if } \mathbf{x} \text{ is the all-one string } 1 \cdots 1 \text{ or the all-zero string } 0 \cdots 0, \\ 2n - 2 & \text{else} \end{cases}$$

where $|\mathbf{x}|$ is the length of \mathbf{x} . For $n \geq 4$, the results in Theorem 4 is stronger than the above one.

Moreover, since all the suffixes of \mathbf{x} including the null string are contained in $\mathcal{D}(\mathbf{x})$, we have $\#\mathcal{D}(\mathbf{x}) \geq n + 1$. Therefore, we obtain the following corollary.

Corollary 2. *For $\mathbf{x} \in \mathcal{A}^n$,*

$$\#\mathcal{AD}(\mathbf{x}) \leq \#\mathcal{D}(\mathbf{x}).$$

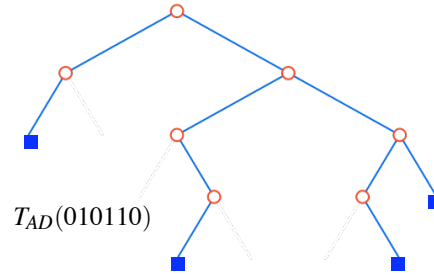


Fig. 3: An example of digital trie $T_{AD}(010110)$.

5 From Antidictionary to Dictionary

The antidictionary $\mathcal{AD}(x)$ of a given string $x \in \mathcal{A}^*$ can be represented as a digital trie denoted by $T_{AD}(x)$. As an example, Figure 3 shows a digital trie representing the antidictionary of $x = 010110$ where $\mathcal{AD}(010110) = \{00, 1010, 1101, 111\}$. An MFW in $\mathcal{AD}(x)$ is corresponding to a path from the root to a leaf in $T_{AD}(x)$.

In Algorithm **AD2D** described below, starting with $T_{AD}(x)$, nodes in $T(x)$ will be reproduced step by step. However, it will be determined on each process of the algorithm whether those nodes are internal or external nodes in $T(x)$. Hence, nodes created by the algorithm are tentatively called ‘neutral’ nodes until we know their connectivity, that is, the number of their child nodes and their directions. Besides, at the initial state of the algorithm, internal nodes in $T_{AD}(x)$ should be treated as neutral nodes since their connectivity are not known at the initial state.

(Algorithm **AD2D**)

0. Let T be $T_{AD}(x)$.
1. For each neutral node p in T in the breadth-first order:
 - 1-1. If $\sigma(p)$ has two internal nodes as child nodes, create one or two new nodes connecting to p so that p has two child nodes.
 - 1-2. If $\sigma(p)$ has only one internal node q as a child node, create a new node connecting to p so that p has a child node with the same direction as q .
 - 1-3. If $\sigma(p)$ has no child nodes, let p be an external node.
2. Remove all the external nodes in T_{AD} from T .
3. Output T as $T(x)$. Then stop.

Figure 4 depicts the process of reconstructing the dictionary of $x = 010110$ from its antidictionary.

Theorem 5. *The suffix trie $T(x)$ is reproduced from $T_{AD}(x)$ by means of Algorithm **AD2D**.*

Proof. Let p be an internal node in $T_{AD}(x)$. From the definition of MFW, $w(p)$ is in $\mathcal{D}(x)$. Suppose that p has only a child node. Then, Algorithm **AD2D** sprouts a new child node q from p if $\sigma(q)$ is not an external node in T . The new node q becomes neutral and $w(q) \in \mathcal{D}(x)$.

The connectivity of a neutral node r is determined by that of $\sigma(r)$. That is, if $\sigma(r)$ has k child nodes in $T(x)$ where $k = 0, 1, 2$, then r does so. Since the algorithm processes neutral nodes in the breadth-first order, the connectivity of $\sigma(r)$ is known when r is processed. Thus, all the suffices of x will be reproduced in Step 1 of the algorithm. Removing all the external nodes in T_{AD} from T , we obtain $T(x)$. \square

Since x is equal to $w(p)$ for a certain leaf p that has the longest path among leaves in $T(x)$, Algorithm **AD2D** can reproduce x from $\mathcal{AD}(x)$. Hence we have the following corollary.

Corollary 3. *The original string x can be reproduced from $\mathcal{AD}(x)$.*

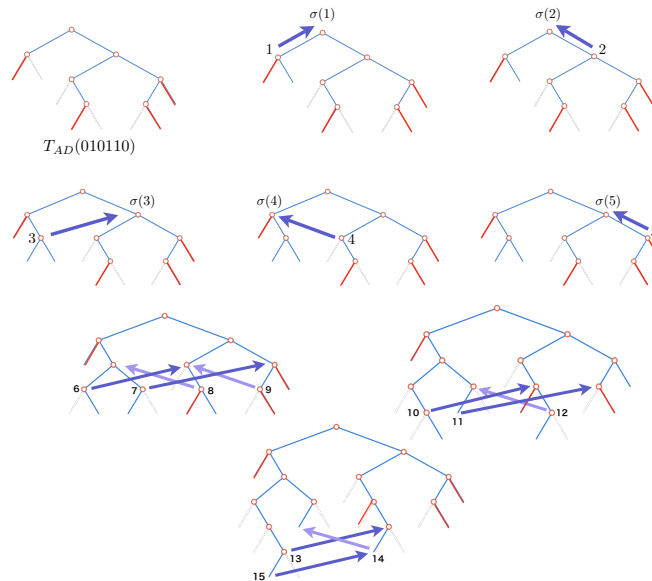


Fig. 4: $T(0101110)$ is reproduced from $T_{AD}(0101110)$ by Algorithm **AD2D** where non-external nodes are indexed by numbers.

6 Conclusions

In this article, we derived an upper bound on the size of the antedictionary of a given binary string x . And we proved that the antedictionary of x is always smaller than or equal to the dictionary of x . Moreover, we gave an algorithm to reconstruct the dictionary of x from the antedictionary of x .

Acknowledgements

The authors express their thanks to Philippe Flajolet of INRIA, Rocquencourt, France to call their attention to Janson et al. (2004).

References

- M. Crochemore, F. Mignosi, and A. Restivo. Automata and forbidden words. *Information Processing Letters*, 67(3):111–117, 1998.
- M. Crochemore, F. Mignosi, A. Restivo, and S. Salemi. ata compression using antedictionaries. *Proc. IEEE*, 88(11):1756–1768, 2000.
- D. Gusfield. *Algorithms on strings, trees, and sequences: Computer Science and Computational Biology*. Cambridge Univ. Press, 1997.
- S. Janson, S. Lonardi, and W. Szpankowski. On average sequence complexity. *Theoretical Computer Science*, 326:213–227, 2004.
- T. Ota and H. Morita. One-path ecg lossless compression using antedictionaries. *IEICE Trans. Fundamentals (Japanese Edition)*, J87-A(9):1187–1195, 2004.