

On Antidictionary Coding Based on Compacted Substring Automaton

Takahiro Ota

Dept. of Computer & Systems Engineering
Nagano Prefectural Institute of Technology
Ueda, Nagano 386-1211, Japan
Email: ota@pit-nagano.ac.jp

Hiroyoshi Morita

Grad. School of Information Systems
University of Electro-Communications
Chofu, Tokyo 182-8585, Japan
Email: morita@is.uec.ac.jp

Abstract—Lossless data compression via substring enumeration (CSE) has been proposed by Dubé and Beaudoin in 2010. The CSE outputs its encoder called compacted substring automaton as a codeword, and an efficient representation of the automaton is also proposed. In this paper, we prove an isomorphism between compacted substring automaton and antidictionary automaton, which is an encoder of antidictionary coding. Then we propose a new static antidictionary coding which uses the representation of compacted substring automaton instead of antidictionary. Moreover, we prove an asymptotic optimality of the proposed coding for a stationary ergodic source.

I. INTRODUCTION

An antidictionary coding is a lossless data compression using an antidictionary automaton which accepts all strings that contain no strings of antidictionary as their substrings. An antidictionary is a set of minimal forbidden words of an input string, and an antidictionary automaton was originally introduced by Crochemore *et al.* in [1]. Various antidictionary codes have been proposed [2]–[4], and the asymptotic optimality for a proper Markov source has been proved [5]. Lossless data compression via substring enumeration (CSE) [6] is also a lossless data compression algorithm using a compacted substring automaton that accepts all substrings of circular string of an input binary string. The compression algorithm constructs a compacted automaton representation of a binary input string as its codeword. The asymptotic optimality for sequences satisfying a proper condition on stationary ergodic source has been proved [7].

In a static antidictionary coding, an antidictionary of a given input string is output as a codeword. An efficient representation of the antidictionary by means of a recursive tree representation has been proposed [2], however, the cost is too large to represent the antidictionary in practice.

In this paper, we first prove an isomorphism between both the encoders that are compacted substring automaton and antidictionary automaton. Then we propose a new antidictionary code which outputs an efficient representation of compacted substring automaton as its codeword instead of antidictionary. Moreover, we evaluate probability of a set of sequences not-satisfying the proper condition shown in [6] for the CSE algorithm on a stationary ergodic source. We prove that asymptotic optimality of the proposed algorithm for a stationary ergodic source.

II. BASIC NOTATIONS AND DEFINITIONS

Let \mathcal{X} be a binary source alphabet $\{0, 1\}$. Let \mathcal{X}^* be the set of all finite strings over \mathcal{X} , including the empty string of length zero, denoted by λ , and let $\mathcal{X}^+ = \mathcal{X}^* \setminus \{\lambda\}$. For convenience, we define $\bar{a} = 1 - a$ for $a \in \mathcal{X}$.

The length of a string x is denoted as $|x|$. We also use $|\cdot|$ to represent the cardinality of a set. For a given string x , by letting $n = |x|$, a substring x_i^j is defined as

$$x_i^j = \begin{cases} x_i \cdots x_j & (1 \leq i \leq j \leq n), \\ \lambda & (i > j). \end{cases} \quad (1)$$

Hereinafter, with no notice, we assume that the length of x is always given as n and $n \geq 2$. Therefore, the substring x_1^n of x equals x . A dictionary $\mathcal{D}(x)$ is defined as the set of all the substrings of x , that is,

$$\mathcal{D}(x) = \{x_i^j | 1 \leq i \leq j \leq n\} \cup \{\lambda\}. \quad (2)$$

Let $\mathcal{P}(x)$ and $\mathcal{S}(x)$ respectively denote the set of all the prefixes and the suffixes of x .

$$\mathcal{P}(x) = \{x_1^i | 1 \leq i \leq n\} \cup \{\lambda\}, \quad (3)$$

$$\mathcal{S}(x) = \{x_j^n | 1 \leq j \leq n\} \cup \{\lambda\}. \quad (4)$$

For a given x_1^n , let $\pi(x_1^n)$ and $\sigma(x_1^n)$ respectively denote $x_1^{n-1} \in \mathcal{P}(x_1^n)$ and $x_2^n \in \mathcal{S}(x_1^n)$. For convenience, we define $\pi(\lambda) = \sigma(\lambda) = \lambda$. For a non-negative integer k , let $\pi^k(x_1^n)$ and $\sigma^k(x_1^n)$ be the functions $\pi(\cdot)$ and $\sigma(\cdot)$ applied k times to a string x_1^n , where $\pi^0(x_1^n) = \sigma^0(x_1^n) = x_1^n$.

A. Necklace

For a given string $x \in \mathcal{X}^+$ and $k \geq 1$, let $x^{(k)}$ be a string that is generated by concatenating x followed by $x^{(k-1)}$ where $x^{(0)} = \lambda$. As an example, for $k = 2$, $x^{(2)} = xx = x_1 \dots x_n x_1 \dots x_n$. A necklace of x is defined as $x^{(\infty)}$. In other words, the necklace of x is a circular string for x .

We give Condition 1 for x . Hereinafter, with no notice, we assume that x satisfies Condition 1 in the following discussions except Section V.

Condition 1. $|\{u \in \mathcal{D}(\sigma(x^{(2)})) | n = |u|\}| = n$.

In other words, for $1 \leq i \neq j \leq n$, a substring of length n starting from i in $x^{(2)}$ is not equal to a string of length n starting from j in $x^{(2)}$. Moreover, x that satisfies Condition 1

Let \mathcal{B}_n be the set of all the strings that satisfy Condition 1 in \mathcal{X}^n . For $x \in \mathcal{B}_n$,

$$\mathcal{W}_n(\mathbf{x}) = \{\mathbf{u} \mid \|\mathbf{u}\| = n, \mathbf{u} \in \mathcal{D}(\sigma(\mathbf{x}^{\langle 2 \rangle}))\}. \quad (5)$$

The size of $\mathcal{W}_n(\mathbf{x})$ is equal to n since \mathbf{x} satisfies Condition 1. The set of all the substrings of length n of $\mathbf{x}^{(\infty)}$ is equal to $\mathcal{W}_n(\mathbf{x})$.

B. Antidictionary

A string v_1^k ($k \geq 1$) $\in \mathcal{X}^*$ with the following three properties

$$\mathbf{v}_1^k \notin \mathcal{D}(\mathbf{x}) \quad (6)$$

$$\mathbf{v}_1^{k-1} \in \mathcal{D}(\mathbf{x}) \quad (7)$$

$$\mathbf{v}_2^k \in \mathcal{D}(\mathbf{x}) \quad (8)$$

is called a *Minimal Forbidden Word (MFW)* of x . An antidictionary $\mathcal{A}(x)$ is the set of all the MFWs of x . For convenience, we define $\mathcal{D}(\lambda) = \{\lambda\}$ and $\mathcal{A}(\lambda) = \mathcal{X}$. Let $\mathcal{A}_m(x)$ be a subset of $\mathcal{A}(x)$ in which the length of an element is shorter than or equal to m . For example, $\mathcal{A}(x)$ for $x = 01000001$ is given as

$$\mathcal{A}(\mathbf{x}) = \{11, 101, 0010, 1001, 10001, 000000, 100001\}. \quad (9)$$

Moreover, for $\mathbf{x}^{(2)} = 0100000101000001$,

$$\mathcal{A}_8(\mathbf{x}^{(2)}) = \{11, 1001, 00100, 10001, 10101, 000000, 100001\}. \quad (10)$$

C. Core

For $\mathbf{v} = a\mathbf{w}b \in \mathcal{X}^+$ and $a, b \in \mathcal{X}$, $c(\mathbf{v}) = \mathbf{w} \in \mathcal{X}^*$ is called *core* of \mathbf{v} . For $\mathbf{t} \in \mathcal{D}(\mathbf{x}^{(2)})$ and $\mathbf{u} = c(\mathbf{t})$, the set of all \mathbf{u} with the following properties

$$0u, 1u, u0, u1 \in \mathcal{D}(x^{\langle 2 \rangle}) \quad (11)$$

is represented by $\mathcal{C}(x)$. For example, $x = 01000001$,

$$\mathcal{C}(\mathbf{x}) = \{\lambda, 0, 00, 000, 010, 0000\}. \quad (12)$$

D. Antidictionary Tree and Automaton

For a given $\mathcal{A}_n(\mathbf{x}^{(2)})$, an *antidictionary tree (AD-tree)* $\mathbb{T}_A(\mathbf{x}^{(2)})$ is a tree structure that stores all the elements of $\mathcal{A}_n(\mathbf{x}^{(2)})$ [1]. Fig. 1 shows $\mathbb{T}_A(\mathbf{x}^{(2)})$ for $\mathbf{x} = 01000001$, where $\mathcal{A}_8(\mathbf{x}^{(2)})$ is shown in (10). In Fig. 1, solid circles and

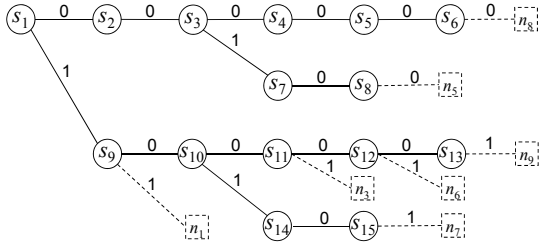


Fig. 1. AD-tree $\mathbb{T}_A(\mathbf{x}^{\langle 2 \rangle})$ for $\mathbf{x} = 01000001$.

dotted squares respectively represent internal nodes and external nodes called *leaves*. The solid and dotted lines respectively

represent edges between internal nodes and between internal nodes and leaves.

In a tree structure, a string associated with a path from the root ρ to a node p in the tree is denoted as $w(p)$. Then we define that $w(\rho)$ is λ . The string $w(p)$ is called *path-string* of p . On the other hand, for a string \mathbf{u} , let $l(\mathbf{u})$ be the node p such as $w(p) = \mathbf{u}$, and $l(\mathbf{u})$ be called *locus* of \mathbf{u} . For example, in Fig. 1, $w(s_8)$ is 0010 and $l(1010)$ is s_{15} ; s_1 is the root ρ .

An *antidictionary automaton* (*AD-automaton*) $\mathbb{A}(\mathbf{x}^{(2)})$ is a deterministic automaton constructed from $\mathbb{T}_A(\mathbf{x}^{(2)})$ by means of the L-automaton algorithm shown in [1], [2]. An AD-automaton $\mathbb{A}(\mathbf{x}^{(2)})$ accepts all strings that contain no strings of $\mathcal{A}_n(\mathbf{x}^{(2)})$ as their substrings.

A state of $\mathbb{A}(\mathbf{x}^{(2)})$, corresponding to an internal node of $\mathbb{T}_A(\mathbf{x}^{(2)})$, has two outgoing edges. On the other hand, a state corresponding to a leaf, called *sink*, has no outgoing edge. Edges are defined in the following manner: for each internal node p and $a \in \mathcal{X}$,

- (i) if $l(\mathbf{w}(p)a)$ exists in $\mathbb{T}_A(\mathbf{x}^{(2)})$, then the edge labeled a from $l(\mathbf{w}(p))$ terminates at $l(\mathbf{w}(p)a)$.
- (ii) if $l(\mathbf{w}(p)a)$ does not exist in $\mathbb{T}_A(\mathbf{x}^{(2)})$, then the edge labeled a from $l(\mathbf{w}(p))$ terminates at $l(\mathbf{v})$, where \mathbf{v} is the longest suffix of $\mathbf{w}(p)a$ and $\mathbf{v} = \mathbf{w}(q)$ with node q of $\mathbb{T}_A(\mathbf{x}^{(2)})$.

Fig. 2 shows $\mathbb{A}(x^{(2)})$ for $x = 01000001$, where $\mathcal{A}_8(x^{(2)})$ is shown in (10). In Fig. 2, dotted triangles and edges respec-

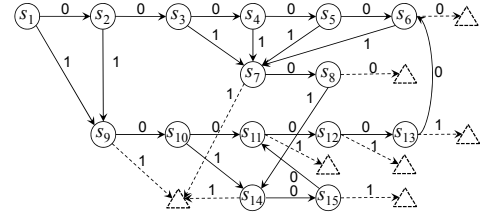


Fig. 2. AD-automaton $\mathbb{A}(\mathbf{x}^{(2)})$ for $\mathbf{x} = 01000001$.

tively represent sinks and edges to sinks.

E. Compacted Substring Tree and Automaton

For a given $\mathcal{D}(\mathbf{x}^{(2)})$, a *compacted substring tree (CS-tree)* $\mathbb{T}_C(\mathbf{x})$ is a tree structure that stores all the strings $0\mathbf{u}$ and $1\mathbf{u}$ such that both $0\mathbf{u} \in \mathcal{D}(\mathbf{x}^{(2)})$ and $1\mathbf{u} \in \mathcal{D}(\mathbf{x}^{(2)})$. Fig. 3 shows $\mathbb{T}_C(\mathbf{x})$ for $\mathbf{x} = 01000001$.

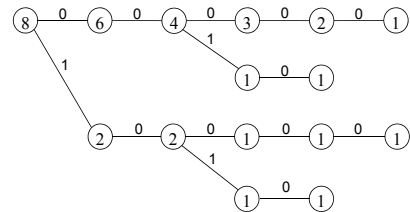


Fig. 3. CS-tree $\mathbb{T}_C(x)$ for $x = 01000001$.

For a node p of $\mathbb{T}_C(\mathbf{x})$,

$$N(\mathbf{w}(p)) = |\{\mathbf{u} \mid \mathbf{w}(p) \in \mathcal{P}(\mathbf{u}), \mathbf{u} \in \mathcal{W}_n(x)\}|. \quad (13)$$

In Fig. 1, a number written in a node p represents $N(w(p))$.
As for $\mathbb{T}_C(x)$, Theorem A and Theorem B hold [9].

Theorem A (Theorem 1 [9]). *For $\mathbb{T}_C(x)$, there exists a node $l(0u)$ if and only if $0u, 1u \in \mathcal{D}(x^{(2)})$. Similarly, there exists a node $l(1u)$ if and only if $0u, 1u \in \mathcal{D}(x^{(2)})$.*

Theorem B (Theorem 2 [9]). *For $\mathbb{T}_C(x)$, a subtree having $l(0)$ as the root and a subtree having $l(1)$ as the root are isomorphic, both of which have $n - 1$ nodes.*

A compacted substrating automaton (CS-automaton) $\mathbb{C}(x)$ is a deterministic automaton that is constructed from $\mathbb{T}_C(x)$. A CS-automaton $\mathbb{C}(x)$ accepts all the strings in $\mathcal{D}(x^{(\infty)})$. A CS-automaton was called CST in an earlier report [6].

A state of $\mathbb{C}(x)$, corresponding to a node of $\mathbb{T}_C(x)$, has one or two edges. Edges are defined in the following manner: for each node p and $a \in \mathcal{X}$,

- (i) if $l(w(p)a)$ exists in $\mathbb{T}_C(x)$, then the edge labeled a from $l(w(p))$ terminates at $l(w(p)a)$.
- (ii) if $l(w(p)a)$ does not exist in $\mathbb{T}_C(x)$ and $w(p)a \in \mathcal{D}(x^{(2)})$, then the edge labeled a from $l(w(p))$ terminates at $l(u)$, where u is the longest suffix of $w(p)a$ and $u = w(q)$ with node q of $\mathbb{T}_C(x)$.

In [6], an edge of (i) is called *forward edge* and an edge of (ii) is called *backward edge*. Fig. 4 shows $\mathbb{C}(x)$ for $x = 01000001$.

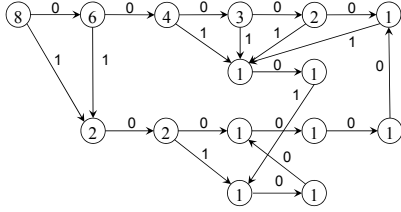


Fig. 4. CS-automaton $\mathbb{C}(x)$ for $x = 01000001$.

III. REVIEWS OF ANTIDictionary AND CSE CODINGS

A. Antidictionary Coding

In a static antidictionary coding, an AD-automaton is used as an encoder, and we use $\mathbb{A}(x^{(2)})$ to compress an input string x in this paper. From the initial state s_1 , the transitions with x are implemented. In the transitions, a symbol is output in a state having two outgoing edges to non-sinks, while no symbol is output in a state having one outgoing edge to a sink. In a static encoding, a transition to sink never occurs since the transition corresponds to an occurrence of an MFW. Therefore, the transition from a state having only one outgoing edge to non-sink is predictable.

The static antidictionary coding outputs the following triples

$$(n, \mathcal{A}_n(x^{(2)}), \gamma(x)) \quad (14)$$

where $\gamma(x)$ represents a string concatenating output symbols [2]. As an example, for $x = 01000001$ and $\mathbb{A}(x^{(2)})$ shown in Fig. 2, $(8, \mathcal{A}_n(x^{(2)}), 010)$ is output.

B. CSE Coding

In the CSE coding, a CS-automaton is used as an encoder. In the decoding, $\mathbb{C}(x)$ is constructed from $\Phi(x)$,

$$\Phi(x) = \phi(N(0))\phi(N(0w_10)) \dots \phi(N(0w_m0)) \quad (15)$$

where $\mathcal{C}(x) = \{w_1, \dots, w_m\}$ and $\phi(i)$ is a binary representation of i . The CSE outputs the following triples

$$(n, \Phi(x), \text{rank}(x)) \quad (16)$$

where $\text{rank}(x)$ represents the rank of x in the lexicographical order of x in $\mathcal{W}_n(x)$ [6]. Note that the set of all the output strings of length n from $\mathbb{C}(x)$ is equal to $\mathcal{W}_n(x)$. As an example, for $x = 01000001$, (12), and Fig. 3, $(8, \phi(6)\phi(4)\phi(3)\phi(2)\phi(1)\phi(0)\phi(0), 4)$ is output.

IV. RELATIONSHIP BETWEEN AD-AUTOMATON AND CS-AUTOMATON

In this section, we will prove the following theorem.

Theorem 1. *For a given x , $\mathbb{A}(x^{(2)})$, without all the sinks and the edges to the sinks, and $\mathbb{C}(x)$ are isomorphic.*

From the constructions of an AD-automaton and of a CS-automaton, Theorem 1 is a corollary of Proposition 1.

Proposition 1. *For a given x , $\mathbb{T}_A(x^{(2)})$, without all the leaves and the edges to the leaves, and $\mathbb{T}_C(x)$ are isomorphic.*

To prove Proposition 1, we first give four lemmas.

Lemma 1. *For $w \in \mathcal{C}(x)$, $|w| \leq n - 2$.*

Proof: For $w \in \mathcal{C}(x)$, $0w \in \mathcal{D}(x^{(2)})$ and $1w \in \mathcal{D}(x^{(2)})$ hold. From Theorem A, both $0w$ and $1w$ are path-strings in $\mathbb{T}_C(x)$. From Theorem B, $|0w| \leq n - 1$ (resp. $|1w| \leq n - 1$) since there are just $n - 2$ edges in a subtree whose root is $l(0)$ (resp. $l(1)$). Therefore, $|w| \leq n - 2$. ■

Lemma 2. *For $u \in \mathcal{A}_n(x^{(2)})$, $c(u) \in \mathcal{C}(x)$.*

Proof: Since $n \geq 2$ and Condition 1 hold, $2 \leq |u| \leq n$ holds. Therefore, for $a, b \in \mathcal{X}$ and $w \in \mathcal{X}^*$, u is written as awb . From (7) and (8),

$$aw, wb \in \mathcal{D}(x^{(2)}). \quad (17)$$

Since $|wb| \leq n - 1$, there exists symbol c such that $cwb \in \mathcal{D}(x^{(2)})$. However, since $awb \notin \mathcal{D}(x^{(2)})$, $cwb = \bar{a}wb$ holds. Therefore,

$$\bar{a}w \in \mathcal{D}(x^{(2)}). \quad (18)$$

Similarly, since $|aw| \leq n - 1$, there exists symbol d such that $awd \in \mathcal{D}(x^{(2)})$. However, since $awb \notin \mathcal{D}(x^{(2)})$, $awd = aw\bar{b}$ holds. Therefore,

$$w\bar{b} \in \mathcal{D}(x^{(2)}). \quad (19)$$

From (17), (18), and (19), $w = c(u) \in \mathcal{C}(x)$. ■

Lemma 3. *For $w \in \mathcal{C}(x)$, there exists $a \in \mathcal{X}$ and $u, v \in \mathcal{X}^+$ such that $awu \in \mathcal{A}_n(x^{(2)})$ and $\bar{a}wv \in \mathcal{A}_n(x^{(2)})$.*

Proof: From Lemma 1, $|w| \leq n - 2$. Since $w \in \mathcal{C}(x)$, there exists $t \in \mathcal{X}^+$ such that $awt \in \mathcal{D}(x^{(2)})$ and $|awt| =$

n . From Theorem B, $\bar{a}wt \notin \mathcal{D}(x^{(2)})$ since $|wt| = n - 1$. Therefore, there exists $k \geq 0$ such that $\pi^k(\bar{a}wt) \notin \mathcal{D}(x^{(2)})$ and $\pi^{k+1}(\bar{a}wt) \in \mathcal{D}(x^{(2)})$ since $\bar{a}w \in \mathcal{D}(x^{(2)})$. Let v be $\pi^k(t)$. Since $awv \in \mathcal{D}(x^{(2)})$,

$$\bar{a}wv \notin \mathcal{D}(x^{(2)}) \quad (20)$$

$$\pi(\bar{a}wv) \in \mathcal{D}(x^{(2)}) \quad (21)$$

$$\sigma(\bar{a}wv) = wv \in \mathcal{D}(x^{(2)}) \quad (22)$$

hold. From (21) and Theorem B, $|\pi(\bar{a}wv)| = |wv| \leq n - 1$. Since (20), (21), (22), and $|wv| \leq n - 1$, $\bar{a}wv \in \mathcal{A}_n(x^{(2)})$.

Similarly, there exists $s \in \mathcal{X}^+$ such that $\bar{a}ws \in \mathcal{D}(x^{(2)})$ and $|\bar{a}ws| = n$. From Theorem B, $aws \notin \mathcal{D}(x^{(2)})$ since $|ws| = n - 1$. Therefore, there exists $l \geq 0$ such that $\pi^l(aws) \notin \mathcal{D}(x^{(2)})$ and $\pi^{l+1}(aws) \in \mathcal{D}(x^{(2)})$ since $aw \in \mathcal{D}(x^{(2)})$. Letting u be $\pi^l(s)$, then since $\bar{a}wu \in \mathcal{D}(x^{(2)})$,

$$awu \notin \mathcal{D}(x^{(2)}) \quad (23)$$

$$\pi(awu) \in \mathcal{D}(x^{(2)}) \quad (24)$$

$$\sigma(awu) = wu \in \mathcal{D}(x^{(2)}) \quad (25)$$

hold. From (24) and Theorem B, $|\pi(\bar{a}wu)| = |wu| \leq n - 1$. Since (23), (24), (25), and $|wu| \leq n - 1$, $\bar{a}wu \in \mathcal{A}_n(x^{(2)})$. ■

Lemma 4. Both $0u \in \mathcal{D}(x^{(2)})$ and $1u \in \mathcal{D}(x^{(2)})$ hold if and only if $u \in \mathcal{P}(v)$ where $v \in \mathcal{C}(x)$.

Proof: For $v \in \mathcal{C}(x)$, $0v, 1v \in \mathcal{D}(x^{(2)})$. Therefore, for $u \in \mathcal{P}(v)$, $0u, 1u \in \mathcal{D}(x^{(2)})$.

For $0u, 1u \in \mathcal{D}(x^{(2)})$, from Theorem B, there exists $z \in \mathcal{X}^*$ and $a \in \mathcal{X}$ such that $0uza, 1uz\bar{a} \in \mathcal{D}(x^{(2)})$ and $|uz| \leq n - 2$. Let v be uz , and $v = uz \in \mathcal{C}(x)$ since $0uz, 1uz, uz0, uz1 \in \mathcal{D}(x^{(2)})$. ■

(*Proof of Proposition 1*): From Theorem A and Lemma 4, the set of all the path-strings of $\mathbb{T}_C(x)$ is given as

$$\{u \mid u \in \mathcal{P}(av), v \in \mathcal{C}(x), a \in \mathcal{X}\}. \quad (26)$$

A set of all path-strings of $\mathbb{T}_A(x)$ without all leaves and edges to leaves is given as

$$\{w \mid w \in \mathcal{P}(az), azb \in \mathcal{A}_n(x^{(2)}), b \in \mathcal{X}\}. \quad (27)$$

For any node p of $\mathbb{T}_C(x)$, from Lemma 3 and Lemma 4,

$$w(p) \in \{w \mid w \in \mathcal{P}(az), azb \in \mathcal{A}_n(x^{(2)}), b \in \mathcal{X}\}. \quad (28)$$

From (26) and (28),

$$\begin{aligned} & \{u \mid u \in \mathcal{P}(av), v \in \mathcal{C}(x), a \in \mathcal{X}\} \\ & \subset \{w \mid w \in \mathcal{P}(az), azb \in \mathcal{A}_n(x^{(2)}), b \in \mathcal{X}\}. \end{aligned} \quad (29)$$

For any node q of $\mathbb{T}_A(x^{(2)})$ without the leaves and the edges to the leaves, from Lemma 2 and Lemma 4,

$$w(q) \in \{u \mid u \in \mathcal{P}(av), v \in \mathcal{C}(x), a \in \mathcal{X}\}. \quad (30)$$

From (27) and (30),

$$\begin{aligned} & \{w \mid w \in \mathcal{P}(az), azb \in \mathcal{A}_n(x^{(2)}), b \in \mathcal{X}\} \\ & \subset \{u \mid u \in \mathcal{P}(av), v \in \mathcal{C}(x), a \in \mathcal{X}\}. \end{aligned} \quad (31)$$

From (31) and (29),

$$\begin{aligned} & \{u \mid u \in \mathcal{P}(av), v \in \mathcal{C}(x), a \in \mathcal{X}\} \\ & = \{w \mid w \in \mathcal{P}(az), azb \in \mathcal{A}_n(x^{(2)}), b \in \mathcal{X}\}. \end{aligned} \quad (32)$$

From (32), the set of all the path-strings of $\mathbb{T}_C(x)$ is the same as the set of all the path-strings of $\mathbb{T}_A(x^{(2)})$ without all the leaves and the edges to the leaves. ■

From Proposition 1, for a node p of $\mathbb{T}_C(x)$ and a node q of $\mathbb{T}_A(x^{(2)})$ without all the leaves and the edges to the leaves, a bijective function f such as $f : p \rightarrow q$ is definable. For nodes s, t and $a \in \mathcal{X}$ in $\mathbb{A}(x^{(2)})$ and $\mathbb{C}(x)$, such that the edge from s to t is labeled by symbol a , let (s, t) be the symbol a .

(*Proof of Theorem 1*): For the labelling symbol $a = (p, q)$ where p is a not sink and q is a sink, $w(p)a \notin \mathcal{D}(x^{(2)})$ since $w(p)a$ has an MFW in $\mathcal{A}_n(x^{(2)})$ as its suffix. Therefore, $N(w(p)a) = 0$.

From both the constructions of $\mathbb{A}(x^{(2)})$ and of $\mathbb{C}(x)$ and Proposition 1, for p and q of $\mathcal{C}(x)$ such that (p, q) exists, there exists a bijective function f such that $(p, q) = (f(p), f(q))$ where $f(p)$ and $f(q)$ are the nodes of $\mathbb{A}(x^{(2)})$ without all the sinks and the edges to the sinks. ■

V. PROPOSED ALGORITHM

The proposed algorithm uses $\Phi(x)$ to represent $\mathcal{A}_n(x^{(2)})$ in (14). The proposed algorithm outputs

$$(0, n, \Phi(x), \gamma(x)) \quad (x \in \mathcal{B}_n) \quad (33)$$

$$(1, n, c, \Phi(x_1^c), \gamma(x_1^c)) \quad (x \notin \mathcal{B}_n) \quad (34)$$

as a codeword for a given string x , and c is the shortest length such that $x = x_1^c x_1^c \dots x_1^c$ and $c < n$. In other words, $x_1^c \in \mathcal{B}_n$ and c is a divisor of n .

Let \mathbf{X} be a stationary ergodic source, and let $H(\mathbf{X})$ be the entropy rate of \mathbf{X} and $H(\mathbf{X}) < \infty$. Then \mathbf{X}_1^n is a sequence of random variables $X_1 \dots X_n$ on \mathbf{X} . The recurrence-time R_l is defined as $R_l = \min\{i > 0 \mid \mathbf{X}_1^l = \mathbf{X}_{1+i}^l\}$.

Let ℓ_n be the codeword length per symbol of the proposed algorithm for a random string of length n . We will prove Theorem 2.

Theorem 2. For any stationary ergodic source \mathbf{X} , ℓ_n converges to $H(\mathbf{X})$ in probability as $n \rightarrow \infty$.

To prove Theorem 2, we use Lemma 5 below and the following theorem:

Theorem C (Theorem 2 [7]). For any stationary ergodic source \mathbf{X} , $\lim_{n \rightarrow \infty} \frac{L_{CSE}(\mathbf{X}_1^n)}{n} = H(\mathbf{X})$ with probability one where $L_{CSE}(\mathbf{X}_1^n)/n$ is the model entropy per bit of the CSE.

We convince the Theorem C, however, the case $x \notin \mathcal{B}_n$ is not considered appropriately in the original proof. Hence, we first prepare a proposition.

Proposition 2. If the entropy rate $H(\mathbf{X})$ of \mathbf{X} is positive, then,

$$\lim_{n \rightarrow \infty} P(\mathbf{X}_1^n \notin \mathcal{B}_n) = 0.$$

Proof: For $x \in \mathcal{B}_n$, there exists a positive integer $k < n$, such that k is a divisor of n , which satisfies $x_1^k = x_{n-k+1}^n$

and $x_1^{n-k} = x_{k+1}^n$. In other words, $x = x_1^k x_1^k \dots x_1^k$. Since $1 \leq k \leq \lfloor \frac{n}{2} \rfloor$, $\lceil \frac{n}{2} \rceil \leq n - k < n$ holds. Then,

$$P(X_1^n \notin \mathcal{B}_n) \leq P\left(R_{\lceil n/2 \rceil} \leq \left\lfloor \frac{n}{2} \right\rfloor\right) \quad (35)$$

$$\leq P\left(R_{n/2} \leq \frac{n}{2}\right) \quad (36)$$

$$= P\left(\frac{\log_2 R_{n/2}}{n/2} \leq \frac{\log_2 n/2}{n/2}\right) \quad (37)$$

Theorem II.5.1 [10] shows $\lim_{m \rightarrow \infty} \frac{\log_2 R_m}{m} = H(\mathbf{X})$ with probability 1. Hence, for arbitrary $\varepsilon > 0$ and $\delta > 0$, there exists n_δ such that for any $m \geq n_\delta$

$$P\left(\bigcup_{k=m}^{\infty} \left\{ \left| \frac{1}{k} \log_2 R_k - H(\mathbf{X}) \right| \geq \varepsilon \right\}\right) \leq \delta,$$

which implies

$$P\left(\left\{ \left| \frac{1}{m} \log_2 R_m - H(\mathbf{X}) \right| \geq \varepsilon \right\}\right) \leq \delta.$$

Since $H(\mathbf{X}) < \infty$, it holds $\frac{\log_2 n/2}{n/2} \leq H(\mathbf{X}) - \varepsilon$ if ε is small and n is sufficiently large. Therefore, letting $m = n/2$, we have

$$P(X_1^n \notin \mathcal{B}_n) \leq P\left(\frac{\log_2 R_m}{m} \leq H(\mathbf{X}) - \varepsilon\right) \leq \delta.$$

Since δ is arbitrarily small, the statement of the proposition holds. ■

Lemma 5. $|\gamma(x)| \leq |w_m| + 1$ holds where w_m is the longest string in $\{w | w0, w1 \in \mathcal{D}(x^{(\infty)})\}$.

Proof: We assume that $|\gamma(x)| > |w_m| + 1$. There exists $ua \in \mathcal{P}(x)$ such that $|\gamma(u)| = |w_m| + 1$ and $|\gamma(ua)| = |w_m| + 2$ where $a \in \mathcal{X}$. Since the state transition with u from the initial state of $\mathbb{C}(x)$ has two outgoing edges, $ua \in \mathcal{D}(x^{(2)})$. Therefore, $u0, u1 \in \mathcal{D}(x^{(2)})$ so that $|u| > |w_m|$. It contradicts of maximality of w_m . Hence, $|\gamma(x)| \leq |w_m| + 1$. ■

Now we give a proof of Theorem 2.

(Proof of Theorem 2): For \mathbf{X} , we consider $H(\mathbf{X}) > 0$ and $H(\mathbf{X}) = 0$. First, we show the proof in case of $H(\mathbf{X}) > 0$. For any $\delta > 0$,

$$\begin{aligned} P(|\ell_n - H(\mathbf{X})| > \delta) \\ &= P(|\ell_n - H(\mathbf{X})| > \delta | \mathbf{X} \in \mathcal{B}_n) P(\mathbf{X} \in \mathcal{B}_n) \\ &\quad + P(|\ell_n - H(\mathbf{X})| > \delta | \mathbf{X} \notin \mathcal{B}_n) P(\mathbf{X} \notin \mathcal{B}_n) \end{aligned} \quad (38)$$

$$\leq P(|\ell_n - H(\mathbf{X})| > \delta | \mathbf{X} \in \mathcal{B}_n) + P(\mathbf{X} \notin \mathcal{B}_n) \quad (39)$$

From Proposition 2, the second term of right-hand side in (39) converges to 0 from as $n \rightarrow \infty$.

Next, we consider the first term of right-hand side in (39). From Theorem C, for $x \in \mathcal{B}_n$, the codeword length in (16) per symbol of the CSE converges to $H(\mathbf{X})$ with probability one. From (16) and (33), the proposed algorithm outputs a pair $(0, \gamma(x))$ instead of $\text{rank}(x)$ against the codeword of the CSE. Therefore, we will prove $\frac{1+|\gamma(x)|}{n}$ converges to 0 in probability as n goes to infinity, so that Theorem 2 holds in case of $x \in \mathcal{B}_n$.

From Lemma 5, $|\gamma(x)| \leq |w_m| + 1$. Then $|w_m| + 1$ is the height of the suffix tree of $x^{(\infty)}$. Moreover, from Lemma 2 [7], $w_m \in \mathcal{C}(x)$. Hence, $|w_m| + 1$ is also the height of $\mathbb{T}_C(x)$.

From Theorem II.5.4 [10], the average height of $\mathbb{T}_C(\mathbf{X})$ is given by $c_1 \log_2 n$ where c_1 is a positive constant with probability one. Therefore, $|\gamma(x)| = c_1 \log_2 n + 1$ holds with probability one. Then, $\lim_{n \rightarrow \infty} \frac{|\gamma(x)|}{n} = 0$ in probability. Therefore, ℓ_n converges to $H(\mathbf{X})$ as n goes to infinity since Theorem C holds and $\frac{1+|\gamma(x)|}{n}$ converges to 0 in probability as n goes to infinity. Theorem 2 holds in case of $H(\mathbf{X}) > 0$.

We next consider in case of $H(\mathbf{X}) = 0$. For x on \mathbf{X} with $H(\mathbf{X}) = 0$ and sufficient large n , x can be written by a substring of $z^{(\infty)}$ where z is a string of constant length c . Therefore, $|\gamma(x)| \leq c$ since the height of $\mathbb{T}_C(x)$ is shorter than and equal to c . Hence, $\lim_{n \rightarrow \infty} \frac{|\gamma(x)|}{n} = 0$ holds.

Therefore, in case of $x \in \mathcal{B}_n$, ℓ_n converges to 0 ($= H(\mathbf{X})$) as n goes to infinity since Theorem C holds and $\frac{1+|\gamma(x)|}{n}$ converges to 0 as n goes to infinity. As for $x \notin \mathcal{B}_n$, the length of a part of codeword $(1, c, \Phi(x_1^c), \gamma(x_1^c))$ is a positive constant since c is a constant. Then length of binary representation of n is $O(\log n)$. Hence, in case of $x \notin \mathcal{B}_n$, ℓ_n also converges to 0 ($= H(\mathbf{X})$) as n goes to infinity. Therefore, Theorem 2 also holds in case of $H(\mathbf{X}) = 0$. ■

VI. CONCLUSION

For a given binary string satisfying Condition 1, the results presented herein prove that an antidictionary automaton, without all the sinks and the edges to the sinks, and a compacted substring automaton are isomorphic. Then we proposed a new antidictionary coding which uses a representation of compacted substring automaton as antidictionary. Moreover, we proved the asymptotic optimality of the proposed algorithm for a stationary ergodic source.

ACKNOWLEDGMENT

The authors would like to thank Dr. Danny Dubé for helpful discussions with respect to construction of the CS-automaton from an input string. This research is partly supported by Grant-in-Aid for Scientific Research(C):24500110.

REFERENCES

- [1] M. Crochemore, F. Mignosi and A. Restivo, "Automata and forbidden words," *Inform. Processing Lett.*, 67(3), pp.111–117, 1998(8).
- [2] M. Crochemore, F. Mignosi, A. Restivo and S. Salemi, "Data compression using antidictionaries," *Proc. IEEE*, 88(11), pp.1756–1768, 2000(11).
- [3] M. Fiala and J. Holub, "DCA using suffix arrays," *Proc. DCC2008*, pp.516, 2008(5).
- [4] T. Ota and H. Morita, "On the adaptive antidictionary code using minimal forbidden words with constant lengths," *Proc. ISITA2010*, pp.72–77, 2010(10).
- [5] T. Ota and H. Morita, "Asymptotic optimality of antidictionary codes," *Proc. ISIT2010*, pp.101–105, 2010(6).
- [6] D. Dubé and V. Beaudoin, "Lossless data compression via substring enumeration," *Proc. DCC2010*, pp.229–238, 2010(3).
- [7] H. Yokoo, "Asymptotic optimal lossless compression via the CSE technique," *Proc. CCP2011*, pp.11–18, 2011(6).
- [8] M. Lothaire, "Applied combinatorics on words," *Cambridge*, p.4, 2005.
- [9] D. Dubé and H. Yokoo, "The universality and linearity of compression by substring enumeration," *Proc. ISIT2011*, pp.1619–1623, 2011(8).
- [10] P.C. Shields, "The ergodic theory of discrete sample paths," *A.M.S.*, pp.154–159, 1996.