

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/359317864>

Universal Sinhala Library: Language Specific Encryption Platform for Sinhala Language

Conference Paper · February 2022

DOI: 10.1109/ICMRE54455.2022.9734085

CITATIONS

0

READS

53

2 authors, including:



[Dinindu Koliya Harshanath Webadu Wedanage](#)

University of Wollongong

4 PUBLICATIONS 2 CITATIONS

SEE PROFILE

Universal Sinhala Library: Language Specific Encryption Platform for Sinhala Language

Dinindu Koliya Harshanath Webadu Wedanage
Smart Infrastructure Facility
University of Wollongong
 Wollongong, Australia
 dkhww937@uowmail.edu.au

Samantha Thelijjagoda
SLIIT Business School
Sri Lanka Institute of Information Technology
 Malabe, Sri Lanka
 samantha.t@sliit.lk

Abstract—Security has become a significant challenge in the modern world. Number science and mathematics opened up a vast path to work with programming models to create innovative mechanisms which improve text encryption. Universal Sinhala Library is a Sinhala language-specific text encryption platform. Platform architecture has been designed to demonstrate every possible combination of the Sinhala alphabet, including comma, space and period. The hypothetical architecture includes every book that ever has been written in Sinhala, and every book that ever could be, including every poem, every scientific paper and every piece of document in Sinhala. The main goal of the research is to create an encryption mechanism for the Sinhala language. Linear congruential generator and extended euclidean algorithm have been used along with the Hull–Dobell Theorem to outline the backbone of the encryption platform. At present, it contains all possible combinations of Sinhala characters virtually. Sinhala text to be encrypted should be searched in the platform, and it will return the location of that particular text in the virtual library architecture, which is the encrypted text string for the searched Sinhala text. It is helpful to the people who are using or working with the Sinhala language, moreover for sharing and transferring Sinhala context securely.

Keywords—text encryption, Sinhala language, security, language-specific encryption, virtual architecture

I. INTRODUCTION

Security is one of the biggest concerns of any text transferred through the internet. The growing capacity of digital encoding has opened up vast new avenues for archiving and distributing texts in virtual space, prompting many to declare the imminent obsolescence of print media, the book included. An exciting correlate to this situation is the revival of interest in and support for the idea of the universal Sinhala library, a virtual collection of every text in existence, albeit imagined as an immense database of digitized material with online accessibility. The core algorithm has to engage with every Sinhala character in the Sinhala alphabet. There are 60 contemporary Sinhala characters, and it contains 20 vowels letters and 40 consonant letters. The core algorithm demonstrates a virtual architecture named “Universal Sinhala Library”, which includes each possible combination uniquely identified by its location in the virtual architecture. This particular location path identifies as the derived encrypted string.

There are various encryption algorithms widely used for security purposes. Image encryption, video encryption and text encryption can be categorised as subsections under encryption.

Some algorithms are specified for a subsection, while some algorithms could be used exclusively. For example, the RC5 encryption algorithm is a fast symmetric block cipher fitting for hardware and software implementations. [1] Natural language processing-based text encryption mechanisms are a vital approach under text encryption [2]. A novel symmetric text encryption algorithm based on a logistic map has been introduced for high-security test encryption, and it can be used in real-time applications [3]. Text encryption has been employed in data compression, which has a significantly improved compression ratio [4]. Cryptographic approaches such as Elliptic Curve Cryptography have also been used to implement text encryption, allowing users to encrypt or decrypt any script with ASCII values [5]. Another ASCII value-based data encryption mechanism was introduced using symmetric key encryption technique [6]. Most of these encryption mechanisms has been derived from cryptography concepts and has not focused on introducing language-specific text encryption mechanism. The use of virtual architecture in an encryption mechanism has not been a concentrated concept in the domain.

Jorge Luis Borges has been introduced a virtual library in his book named “The Garden of Branching Paths”, Virtual library includes every possible phrase in English that has been drafted and has needed to be drafted [7]. He introduced the virtual library as an endless number of hexagons; each hexagon has two facing doors, one door to enter the hexagon and another door to exit the hexagon, while the other four sides of the hexagon contain bookshelves, five to each side (Fig. 1).

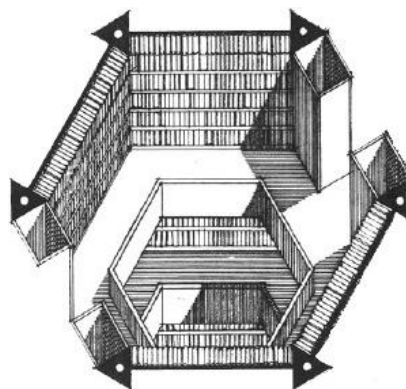


Fig. 1. Hexagon architecture.

Both doors of the hexagon lead to another hexagon. Each hexagon has a hexagon on top of it and a hexagon below it. A spiral staircase located at exiting and entering doors between two adjacent hexagons allows one to walk to the upper level of the library and the lower level of the library [7] (Fig. 2). The library of babel is a website designed by Jonathan Basile. It is a demonstration of Borges' library of babel explained in his book. Currently, it contains all possible pages of 3200 characters, about 10^{4677} books. He has also designed an image archive that contains all possible images that could ever exist [8].

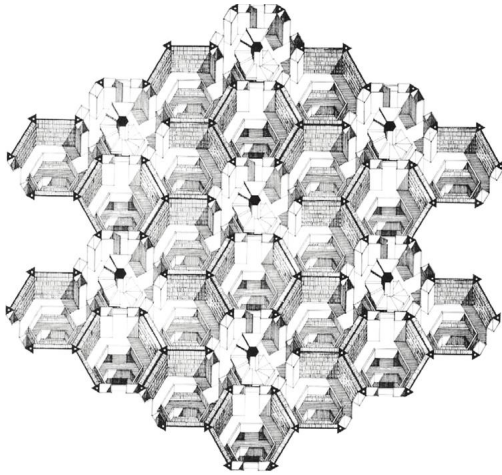


Fig. 2. Floor plan of Borges' Virtual Library, Library of Babel.

Vowels								
අ	ආ	ඇ	ඈ	ඉ	ඊ	උ	ඌ	
a	ā	ae	ae:	i	i:	u	ū	
[a, ə]	[a:, ə]	[æ]	[æ:]	[i]	[i:]	[u]	[u:]	
ඊෂ්ඨ	ඊෂ්ඨ	ඊ	ඊ	ඊ	ඊ	ඊ	ඊ	
ɛ	ɛ	e	e:	ai	o	o:	au	
[ɛ, ɛu]	[ɛ:, ɛu]	[e]	[e:]	[aj]	[o]	[o:]	[əw]	
Vowel diacritics with ka								
ක	කා	කැ	කෑ	කී	කී	කූ	කූ	කා
ka	kā	kae	kaē	ki	kī	ku	kū	ka
කෂ්	කෙ	කේ	කෑ	කො	කෝ	කො	කෝ	කා
kʃ	ke	kē	kai	ko	kō	kau	kañ	kah

Fig. 3. Sinhala vowels and vowel diacritics.

The Sinhala language is derived from the Indo-Aryan group of languages [10]. The Sinhala language has 60 letters, including 20 vowels and 40 consonants (Fig. 4). Apart from that, the Sinhala language has 19 characters which assist in using the Sinhala language. Each consonant comes together with a vowel and forms a vowel diacritic [9] (Fig. 3) Sinhala to English translation demonstrates the Sinhala letters, Singlish letters which means writing Sinhala using English letters and English translation. This demonstration clearly illustrates the formation of vowel diacritics and their use [10].

Universal Sinhala Library is a successful attempt to design a Sinhala language-specific encryption platform inheriting the

Borges' imaginary virtual library and Basile's library of babel website. Primary objectives include design an algorithm to work with Sinhala characters and encrypting Sinhala language content securely. The paper is structured as the introduction in the first section, showcasing the related work and research gaps. Section two demonstrates the overall methodology of the encryption mechanism, the third section illustrates the results, and the last section concludes all sections.

ක	ඛ	ග	ඝ	ඞ	ඟ	ච	ඡ	ජ	ඣ	ඤ	
ka	kha	ga	gha	ṅa	ṅga	ca	cha	ja	jha	ṅa	
[ka]	[ka]	[ga]	[ga]	[ŋa]	[ŋga]	[tʃa]	[tʃa]	[dʒa]	[dʒa]	[ɲa]	
ට	ඨ	ඩ	ඪ	ණ	ඬ	ත	ථ	ද	ධ	න	ඳ
ta	ṭha	ḍa	ḍha	ṇa	ṇḍa	ta	tha	da	dha	na	ṅda
[ta]	[ta]	[ḍa]	[ḍa]	[na]	[ṅḍa]	[ta]	[ta]	[da]	[da]	[na]	[ṅda]
ප	ඵ	බ	භ	ම	ඹ	ය	ර	ල	ව	ළ	
pa	pha	ba	bha	ma	m̐ba	ya	ra	la	va	ḷa	
[pa]	[pa]	[ba]	[ba]	[ma]	[m̐ba]	[ja]	[ra]	[la]	[va]	[la]	
ශ	ෂ	ස	zස	හ	ආ						
śa	ṣa	sa	za	ha	fa						
[ʃa]	[ʃa]	sa	[za]	[ha]	[fa]						

Fig. 4. Sinhala consonants

II. METHODOLOGY

We developed Universal Sinhala Library as a web application where users can search any Sinhala content or upload a softcopy of any Sinhala content for encryption. Initially, it was designed as mentioned in Borges' virtual architecture, and Basile's library of babel [7] [8]. Later on, we made it customizable where users can customize the initial architecture as they prefer. The implementation of the platform is divided into three major components as below.

- Developing base conversion algorithms fitting for the Sinhala language
- Developing the pseudo-random number generation algorithm
- Developing an algorithm to break down Sinhala content into a form where the core algorithm supports.

The above three algorithms work together in the web application; accordingly, all the above algorithms collectively form the core algorithm. The core algorithm functions as a decryption method; therefore, it has a mirror algorithm, which we named the inverted algorithm used as the encryption method.

The Sinhala language has 82 characters, including vowels, consonants, assisting characters, space, comma and period. In order to represent Sinhala content as a number, each character should be given a number; the algorithm is designed to take one page at a time, which means 3200 Sinhala characters; then, any Sinhala content could be represented as a base-82 number. That number is converted into a base-10 number which is fed into the inverted algorithm. The inverted algorithm produces a base-10 unique

output number for the input seed number. That unique base-10 number is broken down into two parts defined in the virtual library architecture: hexagon number and page location. Page location consists of wall number, shelf number, volume number and the page number. Finally, the hexagon number is converted into a base-36 number. The concatenation of the hexagon number and the page location is represented as the encrypted text (Fig. 5).

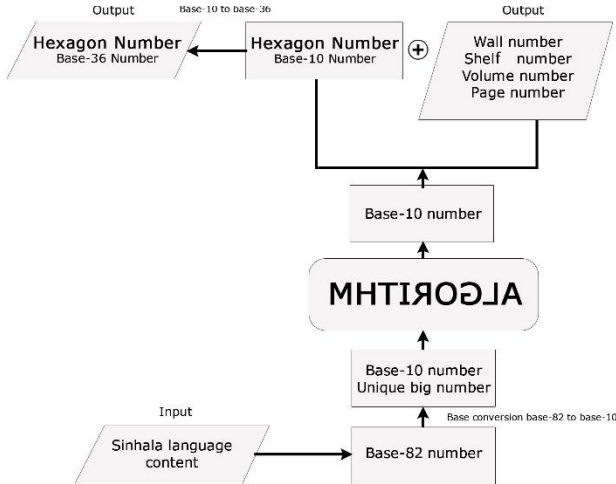


Fig. 5. Encryption algorithm flow chart.

The core algorithm, which is the mirror algorithm of the inverted algorithm, has been employed to decrypt the text. It takes base-36 hexagon number and base-10 page location as inputs. Base-36 hexagon number turned into base-10 number and consolidated base-10 number is taken as the input seed.

The core algorithm produces a unique base-10 number, which is converted into a base-82 number. Then base-82 number is turned into Sinhala content (Fig. 6).

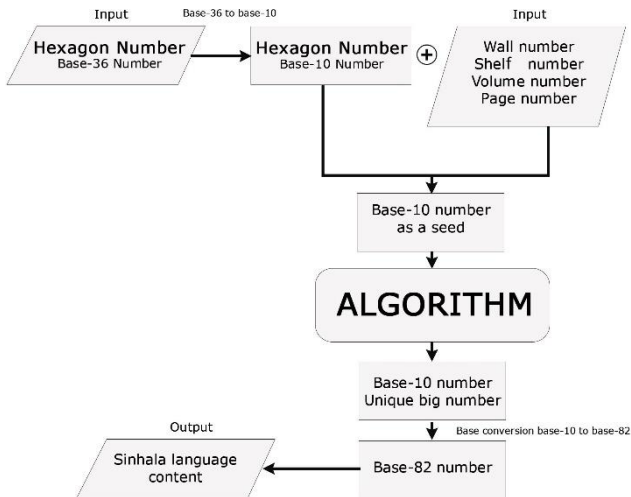


Fig. 6. Decryption algorithm flow chart

The linear congruential generator is used in the core algorithm to yield a sequence of pseudo-randomized numbers determined with a discontinuous piecewise linear equation as shown in equation (1). Recurrence relation has been used to design the pseudo-random number generator [11].

$$X_{n+1} = (aX_n + c) \pmod{m} \quad (1)$$

Where, X is the sequence of pseudorandom values, the modulus (m), $0 < m$, the multiplier (a), $0 < a < m$, the increment (c), $0 < c < m$, and the seed (X_0), $0 < X_0 < m$ are constant integers which defines the pseudorandom number generator.

When $c \neq 0$, accurately selected parameters allow a period equal to m , for all seed values. This will occur if and only if: (1) m and c are relatively prime (2) $a-1$ is divisible by all prime factors of m , and (3) $a-1$ is divisible by 4 if m is divisible by 4. These three requirements are referred to as the Hull–Dobell Theorem.

As shown in equation (2), the extended euclidean algorithm is used to find the accurate coefficients for the above algorithm.

$$ax + by = \gcd(a, b) \quad (2)$$

Extended euclidean algorithm is beneficial when a and b are coprime since x is the modular multiplicative inverse of a modulo b , and y is the modular multiplicative inverse of b modulo a . Extended Euclidean algorithm commonly applied in cryptography [12].

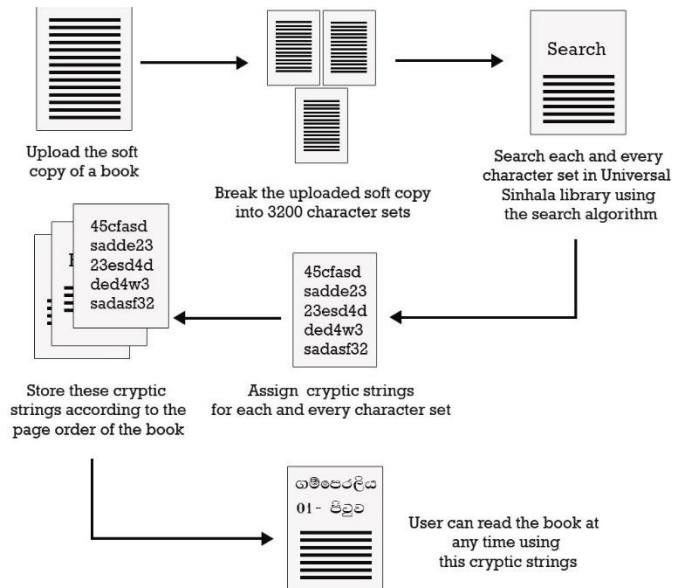


Fig. 7. Breaking down Sinhala content.

Base conversion algorithms and pseudo-random generator algorithms have been demonstrated above, leaving the third component, breaking down Sinhala content into a form that the core algorithm supports (Fig. 7). As mentioned above core algorithm takes only 3200 Sinhala characters once; therefore, Sinhala content is broken down into 3200 Sinhala character strings and merge encrypted strings at the end as the last step of the encryption algorithm. The decryption algorithm is designed to split the encrypted text into encrypted strings, convert each into Sinhala content, and merge them at the end. A user-specific character can be used when merging and splitting the strings.

