

Unicode Searching Algorithm Using Multilevel Binary Tree Applied on Bangla Unicode

Md. Akhtaruzzaman

Department of Computer Science and Engineering, International Islamic University Chittagong Dhaka Campus
akhter.mail@gmail.com, akhter900@yahoo.co.uk

Abstract- Unicode Searching Algorithm using multilevel binary tree is proposed to search the Unicode in efficient way. The algorithm is applied on Bangla Unicode searching to convert Bijoy string into Unicode string. First, the algorithm build a multilevel binary tree based on ASCII code with its corresponding Unicode. The tree is build from a multilevel binary sorted data containing ASCII code and its corresponding Unicode. The data must be sorted based on ASCII code. The algorithm takes Bangla Bijoy string as input value and output the same string in Unicode format. The input Bijoy string must be in Unicode readable format

I. INTRODUCTION

Fundamentally, computers just deal with numbers. They store letters and other characters by assigning a number for each one. Before Unicode was invented, there were hundreds of different encoding systems for assigning these numbers. These encoding systems also conflict with one another. That is, two encodings can use the same number for two different characters, or use different numbers for the same character. Any given computer (especially servers) needs to support many different encodings; yet whenever data is passed between different encodings or platforms, that data always runs the risk of corruption. Unicode is changing these all. Unicode provides a unique number for every character, no matter what the platform, no matter what the program, no matter what the language.

Bangla is the mother tongue of Bangladeshi people and the second most popular language in India. It is a reach language and more than 10% people speak in Bangla in the world. But in field of Computer Science the research on this language is not good enough. There exists some Bangla writing software like Bijoy, Avro, Akkhor etc. Bijoy is the most popular and oldest software which is used to write the plain text only. Avro uses Unicode to write Bangla sentences. Now a day Unicode format is used to write any language and sometimes it is necessary to convert the plain text into Unicode format. So searching Unicode is necessary. In Unicode, there are 65,535 distinct characters that cover all modern languages

of the world. So Unicode searching must be efficient and reliable for all languages. In some Unicode searching algorithm, especially in Bangla, there exist only 'if – else' condition. Some one uses the 'Hash Table' to develop Unicode searching method. Here a Multi Level Tree based Unicode searching algorithm is proposed which will be more efficient for Unicode searching on different languages.

II. PRELIMINARY STUDIES

A. Bijoy String to Unicode Readable Format String

There exist 11 independent characters (vowel) and 39 dependent characters (consonant) [6] in Bangla literature. There also exist some independent and dependent character symbols called 'Kar' and 'Fala' respectively. These symbols must be used with a character. A large number of Complex Characters (combination of two or more characters) exist in Bangla language. A single Complex Character may contain two or more independent or dependent characters, must joined with a symbol named 'Hasanta' (◌̣).

In plain text all the characters and symbols may placed independently anywhere in a sentence. Bijoy follows this rule. But Unicode maintains a unique format to use the symbols with a character. In Unicode the symbols must be used after a character with no gap between them i.e. "character + symbol". But in some cases like 'Chandrabinu' or 'Ref' the placement is different. 'Chandrabinu' must be used after a character if no symbols are exists with that character i.e. "character + Chandrabinu". If any symbol is exist with a character then the 'Chandrabinu' is used after the symbol i.e. "character + symbol + Chandrabinu". 'Ref' must be placed before character.

Figure 1 shows some examples, representing Bijoy plain text and its Unicode readable format of Bangla sentences. Figure 1(a) contains a complex character 'ম্ম' that forms with two ম's, Bangla dependent characters (consonant), i.e. $\text{ম} + \text{্} + \text{ম} = \text{ম্ম}$.

