# A Greedy Approach for Multi-Symbol SMS Message Segmentation

**Liang-Chyau Sheu[1] and Heng Ma[2]**

Department of Industrial Engineering and System Management, Chung Hua University
[1] lcsheu@chu.edu.tw
[2] hengma@chu.edu.tw

## Abstract

Short message service (SMS) has been a very important tool in mobile communication for years, and the usage is still growing. The difficulty in editing messages on handset devices can be overcome by sending SMS messages through Internet-connected gateways. Furthermore, we can adopt such mechanisms to integrate mobile-computing ability to many existing information systems. However, it's easy to send a message in which the length of content exceeds the limit allowed by the SMS specification. Therefore, splitting one message into several segments before passing the message to the SMS gateway is necessary if the original message is lengthy. Due to the specification of SMS, the length of a message segment is determined by the content and may vary in different segments. In this paper, we compare two approaches of segmentation. Proof of optimality for the greedy approach proposed in this paper is an open issue for further research.

**Keywords**: Short Message Service, Segmentation

## 1. Introduction

SMS (Short Message Service) is a technology widely used currently to transmit text message from one mobile phone device to another [1, 2, 3, 4, 5]. Though 3G technology is becoming mature recently, the usage of SMS is still increasing even if it can only carry plain text. There are three major advantages of SMS, which are: (1) SMS is a relatively inexpensive method of communication. (2) SMS is very reliable. (3) SMS is a "written form" communication. As to the first one, though the rates for sending SMS messages vary in different countries and areas, they are significantly lower than those of voice communication in general. For example, the rate for sending one SMS message is NTD$ 3 (some service providers offer their registered customer lower costs), which is far less than that for dialing the number and just saying "Hello!" to your friend. The second advantage is that SMS has been a secure way to send important messages. Even though the recipient's handset is turned off, the message will be stored in the service provider's site and will be sent again later. In most cases, the messages will be delivered to the recipients' handsets within reasonable periods of time. The third advantage can be realized that it is more convenient and practical for some information to be sent and stored in text than voice, e.g. phone numbers and addresses. Because of these advantages, SMS remains competitive among many recently emerged technologies, and will continue to be used in the future.

Conversely, the most significant drawback of sending SMS messages is that editing is relatively difficult on handsets. The numeric key pads used on most mobile phone devices are not easy to input text. Though the youth (commonly referred as "thumb generation" in Taiwan) still are willing to stick on their handsets for sending messages, the inconvenience in editing messages on the handsets has been the major obstacle when using SMS for business purposes. Fortunately, there is an alternative, in which the SMS messages can be sent from Internet-connected computers. Many "gateways," which receive messages in some specific forms (including plain text, URL requests, XML… etc.), have been established. These gateways provide an effective mechanism to embed mobile message dispatching function into existing information systems.

Acknowledging the power of SMS-enabled information system, we have initiated an project in our department (Department of Industrial Engineering & System Management in Chung Hua University) to enhance the instant communication capability in connecting all the people related to the department, including faculty, staff, students and students' parents. So far, we have already integrated SMS services to our web-based campus information system in many functionality fields. We also discovered some issues which require further discussion. For example, the message segmentation problem is critical for long messages, and we will focus on this issue in this paper.

## 2. The Message Segmentation Problem

The term "SMS" reveals that the message length is limited because it's a "short" message. Basically the message length depends on the implementation of the service provider. In general, in Latin-lingual areas, the limitation is 160 characters for one SMS message, which could be insufficient for describing an event in some cases. For those counties and areas that use Chinese in Asia, however, a short message with 70 Chinese characters could be very meaningful. It can be observed that most mobile phone owners in China send tremendous amount of SMS messages. Whereas using computer to send SMS messages, it's very easy to send a long message, because the input methods on computer are much easier to use than the ones on handsets, and there are many editing functions on computer such as "copy/paste". Therefore, we characterize our problem: if there is a lengthy message and segmentation is necessary, we try to determine the optimal segmentation method to send the original message. It seems very easy if the message contains Latin-characters only. However, in Asian countries there is the so-called "two-byte" problem. One Chinese "character" occupies 2 bytes of memory space. It is obvious that mixing Chinese characters with Latin characters in a message will increase the complexity of segmentation. Message segmentation may also be used to solve some other complicated problem, such as exchanging messages of different sizes among a set of processors linked through an interconnected network [6].

The message segmentation problem can be solved by a number of heuristics or algorithms if the length limitation is represented by bytes and the length is fixed. Unfortunately, all the service providers calculate the length of a message with the number of symbols instead of that of bytes. For example, each Latin character is treated as one-byte symbols (Type I), and one Chinese character is also treated as a symbol (Type II) which occupies 2 bytes of memory. Furthermore, the length limitation of a message segment is affected by the types of symbols included in the message and we refer this method as *m/n* scheme hereafter in this paper. For example, 160/70 scheme is commonly used by telecommunication carriers in Taiwan. That is, 160 symbols (referred to as *m*) are allowed in a message segment if they are all Type I symbols. However, if there is any Type II symbol in a message segment, the length limitation becomes 70 (referred to as *n*). This scheme is depicted in Figure 1. We show three situations in Figure 1. However, it's obvious that segment B and C are of the same category.
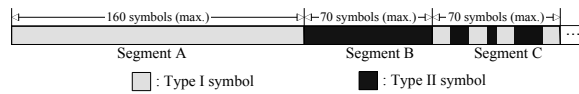


**Fig. 1:** An example of the SMS message segmentation with 160/70 scheme.

With the conditions mentioned above, our objective is to determine the segmentation locations which minimize the number of segments to send the original message.

## 3. The Naïve Approach

With the specification of SMS, it is very easy for someone to implement a message segmentation program in which segments are classified into two types of length. This approach scans the message from the beginning, and the segmentation locations are determined according to a simple rule. The length of the segment is $m$ when there is no Type II symbol within the first $m$ symbols; otherwise, the length of the segment is $n$. One segmentation example is shown in Figure 2 in which 5/2 scheme is adopted for simplicity. The result of this approach is shown in Figure 2(b). It can be noted that a segment with length $n$ might contain symbols which are all of Type I. By observing that, we can conclude that this approach cannot find the optimal solutions in most cases. That's why we call this approach "naïve."
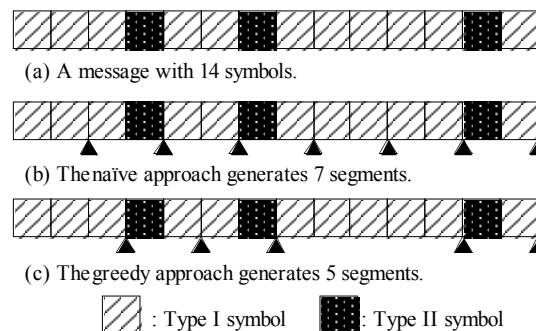


(a) A message with 14 symbols.

(b) The naïve approach generates 7 segments.

(c) The greedy approach generates 5 segments.

**Fig. 2:** Segmentation results of the naïve and greedy approaches with the 5/2 scheme.

## 4. The Greedy Approach

In this approach, the message is scanned form the beginning, and the segmentation locations are determined according to three rules as in the following.

(1) When there is no Type II symbol within the

first $m$ symbols, the length of the segment is $m$, and the next scanning process will start again from the ($m$+1)th symbol.

(2) When the first Type II symbol is found within the first $n$ symbols, the length of the segment is $n$, and the next scanning process will start again from the ($n$+1)th symbol.

(3) When the first Type II symbol is found at position $p$ which is between ($n$+1) and $m$, the length of the segment is ($p$-1), and the next scanning process will start again from the $p$th symbol.

## 5. Conclusion

In this paper, we compared two approaches for dealing with the message segmentation problem. According to the testing examples, we can conclude that the number of segments determined by the greedy approach is always less or equal to the one determined by the naïve approach. However, we could not guarantee the optimality of the greedy approach because it is accomplished by a single scan, and it requires further effort to provide the proof. For this type of problem, the optimal solution can only be obtained by a global inspection of the symbol distribution before making the segmentation.

## 6. References

[1]    Y. Lin, M. Chang, and H. Rao, "Mobile prepaid phone services," *IEEE Personal Communications*, 7(3):6-14, June 2000.

[2]    Jacobsen, K., and Johansen, D., "Ubiquitous Devices United: Enabling Distributed Computing Through Mobile Code," *Proceedings of the Symposium on Applied Computing (ACMSAC'99)*, February 1999.

[3]    S. Guthery, R. Kehr, J. Posegga, and H. Vogt, "GSM SIMs as Web Servers," *Short-Proceedings of 7th International Conference on Intelligence in Services and Networks IS&N'2000*, Athens, Greece, Feb. 2000.

[4]    Nadia Moertiyoso and Kin Choong Yow, "Designing Wireless Enterprise Applications on Mobile Devices," *Proceedings of ICITA2002*. <http://citeseer.csail.mit.edu/554469.html>

[5]    Jouko E. Pakanen, Kai Hakkarainen, Kari Karhukorpi, Petri Jokela, Timo Peltola, and Jyrki Sundström, "A Low-cost Internet Connect for Intelligent Appliances of Buildings," *ITcon* Vol. 7, 2002

[6]    Alfredo Goldman, Joseph G. Peters and Denis Trystram, "Exchanging Messages of Different Sizes," CMPT TR 2002-08, School of Computer Science, Simon Fraser University, September, 2002.