

# 한글 텍스트 스테가노그래피에서 어절을 이용한 정보은닉 기법

## (A Techniques to Conceal Information Using Eojeol in Hanguk Text Steganography)

지 선 수<sup>1)</sup>  
(Ji Seon Su)

**요약** 디지털 시대에 인터넷에서 사용되는 모든 데이터는 디지털화되어 통신 네트워크를 통해 송신 및 수신된다. 따라서 디지털 데이터가 불법적인 사용자에게 의해 변조되고 조작될 수 있기 때문에 기밀성과 무결성을 갖춘 암호화된 데이터를 전송하는 것이 중요하다. 스테가노그래피는 암호화 기법과 혼합되어 기밀성과 무결성을 함께 보장하기 위한 효율적인 방법이다. 커버 매체에 삽입되는 위치와 변화하는 어절 형태를 기반으로 비밀 메시지를 삽입하는 한글 텍스트 스테가노그래피 방법을 제안한다. 한글 텍스트 스테가노그래피에서 3.35%의 삽입용량과 0.4%의 파일 크기 변화를 고려할 때 실험결과는 Jaro\_score 값이 0.946으로 유지할 필요가 있다는 것을 보여준다.

핵심주제어 : 비밀 메시지, 어절크기, 정보은닉, 텍스트 스테가노그래피, 한글 텍스트

**Abstract** In the Digital Age, All Data used in the Internet is Digitized and Transmitted and Received Over a Communications Network. Therefore, it is Important to Transmit Data with Confidentiality and Integrity, Since Digital Data may be Tampered with and Tampered by Illegal users. Steganography is an Efficient Method for Ensuring Confidentiality and Integrity Together with Encryption Techniques. I Propose a Hanguk Steganography Method that Inserts a Secret Message based on a Changing Insertion Position and a Changing Eojeol Size in a Cover Medium. Considering the Insertion Capacity of 3.35% and the File Size Change of 0.4% in Hanguk Text Steganography, Experimental Results Show that the Jaro\_score Value needs to be Maintained at 0.946.

Key Words : Eojeol Size, Hanguk Text, Information Hiding, Secret Message, Text Steganography

### 1. 서론

Corresponding Author : ssji@gwnu.ac.kr

Manuscript received Aug 5, 2017 / revised Sep 6, 2017 / accepted Oct 14, 2017

1) 강릉원주대학교 소프트웨어학과

현대를 살아가면서 필요한 모든 정보는 인터넷을 통해 획득할 수 있으며, 대중의 의견이 수렴되고, 다양한 의혹이 실시간으로 검증되는 도구로 사용되며, 사이버 공간상에서 획기적인 새로운 문화와 현상이 형성되고 있다. 이러한 공간을 기반으로 하는 여러 형태의 사이버 공격 또한 매

우 심각하다. 오늘날 디지털 시대에 온라인상에서 모든 자료는 디지털화되어 통신 네트워크를 통해 송신 및 수신되어진다. 이 과정에서 디지털화된 자료에 불법적으로 접근하거나 해킹 등의 방법으로 위조 및 변조가 이루어질 수 있기 때문에 송신자 및 수신자 모두에게 보안성과 견고성을 강화하기 위한 조치는 매우 중요한 요소이다. 특히 민감 정보를 처리하는 경우에 정보가 훼손되지 아니하도록 안전성 확보에 관련된 조치는 선택이 아니라 필수적인 사항이 되었다.

인터넷 공간에서 자료가 송신 및 수신될 경우 내부적인 자료측면에서 볼 때 보안성이 중요시되며, 외부적인 자료측면에서 무결성이 중요시된다. 텍스트, 이미지, 오디오, 비디오 등과 같이 인터넷상에서 송수신되는 매개체에 은닉된 비밀 메시지를 감지하지 못하게 하여 합법적인 사용자에게 전송하는 스테가노그래피는 암호화 기법과 혼합되어 기밀성과 무결성을 함께 보장하기 위한 효율적인 방법이다. 텍스트 기반 스테가노그래피는 중복 정보의 부족으로 정보를 숨기는 방법이 매우 까다롭지만 인터넷상에 송수신되는 75% 이상의 자료는 텍스트 형태의 자료이기 때문에 가치 있는 연구 분야 중에 하나이다[1-3].

이 논문에서는 한글 어절의 형태에 따라 비밀 문자를 은닉하는 한글 텍스트 스테가노그래피의 변형된 적용 기법을 제시하며, 커버 매체와 스테고 매체 문자열의 유사성을 측정하기 위해 Jaro-Winkler distance를 이용한다.

## 2. 관련 연구

암호화 기법은 텍스트의 패턴과 특성, 은닉여유 공간의 부족 등으로 제3자가 비밀통신을 인식할 수 있는 단점이 있다. 커버 매체의 파일 크기와 자료의 손상이 최소한인 상태를 유지하면서 규모가 비교적 작은 다른 형태의 비밀 정보를 숨기는 스테가노그래피는 암호화 기법과 혼합하여 사용함으로써 보안성과 안정성 모두를 만족하는 효율적인 기법으로 발전되고 있다.

일반적인 텍스트 스테가노그래피는 언어적 스테가노그래피와 기술적인 스테가노그래피로 분류

된다. 언어적인 스테가노그래피는 비밀 메시지를 숨기기 위해 자연언어를 커버 매체로 이용하는 기술이다. 커버 매체를 텍스트로 사용할 경우 코드값이 조금만 변경되어도 의미 자체가 변경되어 지므로 단어의 패턴을 이용하거나 단어 사이의 공백 등에 비트화된 비밀 메시지의 조각을 숨기는 소극적인 방법을 이용하였다. 이러한 방법은 스테고 매체에서 문자의 일부 정보가 수정될 수 있으며, 파일 크기가 증가되는 단점이 있다[3-5].

Monika Agarwal는 하나 이상의 문자가 누락된 단어를 포함하는 접근방법, 단어 목록에 비밀 문자를 숨기는 기법, 암호 메시지를 비트 스트림으로 변환한 후 커버 파일의 단락을 선택하고 숨길 비트에 따라 해당 단어의 시작 또는 끝 문자를 사용하여 각 비트를 숨기는 방법을 제시하였다. 인터넷을 통해 사용자 암호, PIN 등과 작은 크기의 기밀 자료를 안전하게 전송할 때 효율적임을 보였다[3-4]. Kumar는 Monika Agarwal의 은닉 기법을 기반으로 문자에 자료를 은닉하는 방법과 공간에 숨기는 기법을 제안하였다. 비밀 메시지를 비트 스트림으로 변환한 다음에 커버 매체로부터 하나의 문자를 선택한 후 ASCII 값이 짝수 혹은 홀수를 이용하여 은닉될 비트를 생성하는 방법과 ASCII 값이 커버 매체의 공간에 숨겨지며, 공간 조작의 위치를 이용하여 저장될 숫자를 은닉하는 기법을 제시하였다[5-6].

Kingslin과 Kavitha는 Microsoft-word와 Excel 문서에서의 폰트 형태와 문자 순환 기법을 이용하여 텍스트 스테가노그래피에서의 삽입율(embedding ratio/capacity), 불탐지성(undetectability), 견고성(robustness), 불가시성(invisibility)을 평가하였으며, 문자열의 유사성을 측정하였다. MS-word 문서에서 이진화 된 비트를 사용하고 혼합된 폰트 기법을 사용할 경우 유사성은 낮지만 삽입율이 높음을 보였다[7-8]. Monika Agarwal와 Kumar가 제시한 방법에서 삽입하고자 하는 비트화 된 정보에 따라 스테고 키에 정보를 반복적으로 추가해야하는 단점이 있다.

스테가노그래피 기법과 스테간 분석 사이의 끊임없는 도구의 개발 즉, 공격자에 대한 보안적인 대응조치와 취약점 파악의 노력은 향후 계속될

것이다. 따라서 커버 매체에 비교적 작은 크기의 비밀문자 정보를 삽입하는 텍스트 스테가노그래피에서 비밀 메시지의 삽입 전과 후의 내용과 파일 크기의 변동이 최소한으로 유지되는 한글 텍스트 스테가노그래피 적용기법이 필요하다.

이 논문에서는 어절의 크기, 변화하는 삽입 위치, 비밀 메시지의 글자를 초성, 중성, 종성으로 분리한 후 대응되는 이진화 코드를 기반으로 하는 새로운 한글 텍스트 스테가노그래피 기법을 제시한다.

### 3. 분석도구

일반적으로 커버 매체에 삽입되는 정보의량은 삽입 용량과 공간 저장율(saving space ratio) 등을 이용하며, 여기에서는 수식 (1)식을 이용하여 계산한다[3-4].

$$Capacity\ ratio = \frac{\text{숨겨지는 매체 크기(byte)}}{\text{커버 매체의 크기(byte)}} \quad (1)$$

스테고 파일과 커버 파일을 대응시킬 때 일치(matching)된 글자와 전이(transition)가 발생한 글자를 기반으로 하는 유사성(similarity) 측정은 Jaro-Winkler Distance 식을 이용할 수 있다 [3-5][9-10].

$$Jaro - Winkler(C, S) = \frac{Jaro(C, S) + l \cdot p(1 - Jaro(C, S))}{2} \quad (2)$$

혹은

$$Jaro - Winkler(C, S) = \frac{Jaro(C, S) + \frac{\max(p, 4)}{10}(1 - Jaro(C, S))}{2} \quad (3)$$

여기에서 Jaro score는 수식 (4)로 계산할 수 있다.

$$Jaro(C, S) = \frac{1}{3} \left( \frac{m}{|C|} + \frac{m}{|S|} + \frac{m-t}{m} \right) \quad (4)$$

수식에서 사용된  $m$ 은 스테고 문자와 커버 문자에서 매칭된 문자수를 나타내며,  $t$ 는 전이된 문자의 수를 2로 나눈 정수값을 의미한다.  $l$ 은 문자열의 시작에서 최대 4자까지의 공통 접두어의 길이이다.  $p$ 는 공통 접두사를 갖는 점수가 위로 조정되는 정도에 대한 일정한 스케일링 계수이며  $p = 0.25$ 이하로 하며, 보통  $p = 0.1$ 을 이용한다.  $|C|$ 는 커버 매체의 문자열을 나타내고,  $|S|$ 는 스테고 매체의 문자열을 나타낸다.

커버 매체와 비밀 메시지가 삽입된 스테고 매체 사이의 유사성을 비교하기 위한 상관계수는 수식 (5)를 사용하여 계산할 수 있다. 여기에서  $X$ 는 커버 매체,  $Y$ 는 스테고 매체 자료이며,  $x_i \in X$ ,  $y_i \in Y$ 이다.  $\bar{x}$ 와  $\bar{y}$ 는  $X$ 와  $Y$ 의 각각의 평균을 의미한다.

$$Corr = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad (5)$$

### 4. 제안된 방법

한글 문서에서 삽입되는 완성형 문자의 유니코드 정보 대신에 변형된 코드 값을 이용한다.

삽입 위치와 어절의 크기를 이용하며, 초성, 중성, 종성을 기반으로 정보를 삽입하는 기법을 제안한다. 즉 완성형 한글의 유니코드 값 대신에 초성, 중성, 종성 글자로 분리하여 다음의 약속된 코드를 기반으로 각각 4비트 단위의 비트정보를 사용하여 커버 매체에 삽입한다. 한글 자음(초성, 중성)과 모음(중성)에서 사용하는 코드는 Table 1과 같이 나타낼 수 있다. 제안된 알고리즘을 Fig. 1로 나타낼 수 있으며, 그림에서 진하게 표현된 부분이 논문에서 제안된 부분이다. 여기에서 사용되는 커버 매체는 완성형 한글이며, 유니코드 값 U+AC00부터 U+D7A3까지를 고려한다. 또한 전달하고자 하는 비밀 메시지는 문자 개수가 50이하일 경우만을 고려한다.

Table 1 Code used in Hangul consonants (initial, final), vowels(medial)

	Hangul consonants (initial, final)	Hangul vowels (medial)	Binary code
0			0000
1	ㄱ	ㅏ	0001
2	ㅋ	ㅑ	0010
3	ㆁ	ㅓ	0011
4	ㄴ	ㅕ	0100
5	ㄷ	ㅗ	0101
6	ㄹ	ㅛ	0110
7	ㅌ	ㅜ	0111
8	ㅇ	ㅠ	1000
9	ㅈ	ㅡ	1001
10	ㅊ	ㅣ	1010
11	ㅍ	ㅍ	1011
12	ㅍ	ㅍ	1100
13	ㅍ	ㅍ	1101
14	ㅎ	ㅎ	1110
15	-	ㅣ	1111

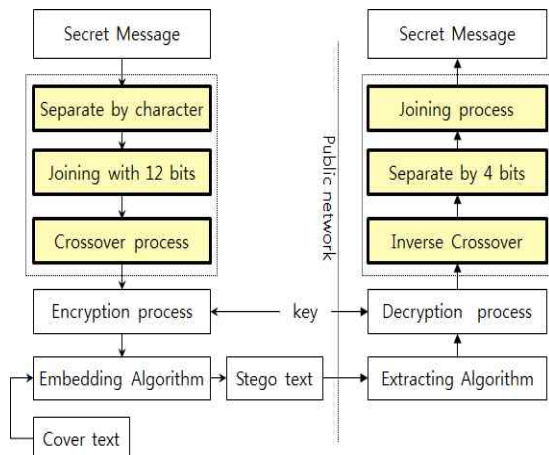


Fig. 1 Text steganography model for the proposed implementation

4.1 은닉 과정

비밀 메시지를 숨기기 위한 은닉 알고리즘은 다음과 같다.

1. 숨기려는 비밀 메시지의 길이를 계산한다. ( $le$ )
2. 비밀 메시지에서 첫 글자를 읽어들이고 다음 유니코드 값을 계산한 후 중간 위치의 값을

초기 은닉시점( $s_o$ )으로 한다.

3. 비밀 메시지로부터 순서대로 문자를 읽어 초성, 중성, 종성 글자로 분리한다. 각각에 대응되는 이진화 코드로 변환한다.
4. 초성, 중성, 종성의 이진화 코드를 결합한 후 앞쪽과 뒤쪽의  $k_1$ 비트를 교차(crossover) 시킨 후 비밀키( $k_2$ )를 이용하여 암호화 한다.
5. 커버 매체에서 임의의 삽입 시점( $s_i$ )을 결정한 후 비밀 메시지의 비트화된 정보를 삽입한다.
  - 5.1 커버 매체에서 비밀 메시지를 삽입할 수 있는지 어절 구조를 확인한다.
  - 5.2 어절의 크기가 홀수이면  $b=0$ 으로, 어절의 크기가 짝수이면  $b=1$ 로 정한다.
  - 5.3 남아있는 비밀 메시지의 비트화된 정보를 삽입하기 위해 5.1부터 5.2의 과정을 12회 반복한다.
  - 5.4  $s_i = 2 * s_{i-1}, i = i + 1$ 를 계산한다.
6. 비밀 메시지의 끝( $le$ )이 나타날 때 까지 3에서 5부분을 반복한다.
7. 스테고키 파일에  $s_o, k_1, k_2, le$ 를 표시한 다음, 수신자의 공개키를 이용하여 암호화한 후( $s_{key}$ ) 스테고 파일과 함께 송신한다.

4.2 추출 과정

스테고 매체로부터 은닉된 비밀 메시지를 찾아내는 추출 알고리즘은 다음과 같다.

1. 수신자의 개인키로  $s_{key}$ 를 복호화 한다. 스테고키 파일로부터  $s_o, k_1, k_2$ 와  $le$ 를 획득한 후  $s_o$ 를 참고하여 은닉시점을 확인한다.
2. 스테고 파일로부터 순서대로 어절을 읽어 들인다. 어절의 크기를 확인하여  $b$ 를 정한다.
3.  $bin[j] = b, j = j + 1$ 를 계산한다.
4. 2부터 3과정을 12회 반복한다.
5. 12개로 이루어진 이진화 코드를 결합한 후 비밀키( $k_2$ )를 이용하여 복호화 한 다음에 앞쪽과 뒤쪽의  $k_1$ 비트를 교차시킨다. 4비트씩

구분하여 이진화 코드 값에 해당되는 초성, 중성, 종성 정보를 가지고 한글 문자를 획득한다.

6.  $bin[j] = 0, j = 0$ 으로 초기화 한다.
7.  $s_i = 2 * s_{i-1}, i = i + 1$ 를 계산한다.
8. 비밀 메시지의 끝( $le$ )이 나타날 때까지 2에서 7부분을 반복한다.

한글 텍스트 정보는 코드값이 조금만 변경되어도 의미 자체가 변경되어지므로 어절의 수, 어절 사이의 공백의 수, 단어의 중성자 유무 등에 비트화된 비밀 메시지의 조각을 숨기는 소극적인 방법을 이용할 수 밖에 없다. 따라서 작은 크기의 비밀 메시지를 암호화와 스테가노그래피 기법을 혼합하여 사용하는 것이 안정성과 견고성 측면에서 효율적이다. 또한 영문과 한글 문서에서 유사한 의미를 표현하기 위해 사용되는 글자의 수는 2.48:1의 비율로 이루어지기 때문에 한글 텍스트 스테가노그래피가 효율적일 수 있으며, 메모리 사용측면에서 두 언어는 비슷하다. 특히 한글문서는 홀수 개 어절과 짝수 개 어절 비율이 0.53:0.47으로 구성되어 있어 좋은 불편성(unbiased)을 가지고 있으므로 어절의 크기를 이용하는 것은 타당하다고 볼 수 있다.

### 5. 적용 및 결과

논문에서 사용된 비밀 메시지의 크기는 2, 6, 10, 14, 30 바이트이며 커버 매체의 크기는 2,266 바이트이다.

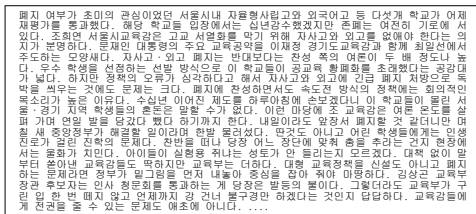


Fig. 2 The cover file

Fig. 2은 커버 매체로 사용한 파일이며, Fig. 3

은 커버 매체에 비밀 문자가 삽입된 스테고 매체 파일이다. 특수 문자를 제외한 순수 한글로 커버 매체와 비밀 문자가 구성되었으며, 알고리즘을 구현하는 과정은 J2SE를 이용하였다.

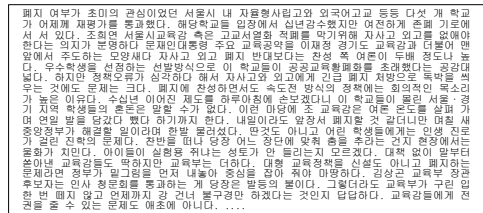


Fig. 3 The stego file

비밀 문자는 ‘대한민국과이팅합시다...’으로 하였으며, 4.1에서 제안된 방법으로 비밀 문자를 은닉하였다. 은닉된 문자를 추출할 경우 4.2의 과정을 이용할 수 있다. 예를 들어 ‘한’을 은닉하기 위해 ‘ㅎ:1110’+‘ㅏ:0001’+‘ㄴ:0010’으로 글자를 분리한 후 커버 매체의 어절 형태를 기반으로 4.1절의 5단계에 따라 이진화된 코드 정보를 은닉한다. 여기에서는 첫 번째 은닉문자가 ‘대(45824)’이므로  $s_0 = 8$ 이며, 임의로  $k_1 = 4, k_2 = 001100110011$ 로 하였다. 커버 매체와 스테고 매체의 유사성을 비교하기 위해 (3)식과 (4)식을 이용하였다.

하나의 비밀 문자를 삽입하기 위해 40개의 어절과 109바이트의 커버 매체가 필요하며, 2.8%의 평균 삽입용량을 유지하며, 평균 4개의 전이가 발생하였으며, 비밀 메시지가 삽입될 경우 2바이트 이하로 파일 크기가 변하는 것이 확인되었다.

Table 2 Results after the embedding process

secret messages	capacity (%)	Jaro_sc ore	Jaro_Winkler	Corr.
1	1.60	0.972	0.978	0.886
3	2.08	0.955	0.964	0.787
5	3.35	0.946	0.956	0.770
7	3.85	0.948	0.958	-
15	4.53	0.921	0.938	-

Table 2에서 커버 매체와 스테고 매체의 유사성과 삽입용량을 확인할 수 있다. 동일한 조건에서 한글 문서 30개를 가지고 각각의 비밀 메시지

가 삽입될 때 비밀 메시지와 Jaro\_score 값은 근사식  $y = -0.00206x + 0.9718$ 으로 표현할 수 있음을 확인하였다. 여기에서  $x$ 는 비밀 메시지의 크기(byte)를 나타낸다. 한글 텍스트 스테가노그래피에서 3.35% 내외의 삽입용량, 0.4%의 파일 크기 변화를 고려할 때 Jaro\_score 값이 0.946으로 유지할 필요가 있으며, 유사도 측면에서 양호함을 확인할 수 있다. 비밀정보 삽입에 따라 스테고 파일에서의 한글 유니코드 값의 변동 폭이 크기 때문에 상관계수는 낮게 나타난다.

## 6. 결 론

비밀 메시지를 완성형 문자의 유니코드를 이용하는 것보다 초성, 중성, 종성으로 분리하여 삽입하면 25%의 저장 공간이 절약될 수 있다. 또한 삽입용량을 3.35% 이하로 유지하여 비밀 문자를 삽입할 경우 Jaro\_Winkler 값이 0.95에 근접하여 유사성 측면에서 효과적임을 확인하였다. 따라서 인터넷에서 송수신되는 대부분의 자료가 텍스트이고, 구조와 형태가 다양하기 때문에 한글 텍스트 스테가노그래피는 보안성이 강화된 작은 용량의 비밀 메시지를 은닉하는데 적절한 도구가 될 수 있다. 어절에서 첫 글자의 초성과 끝 글자의 종성을 이용한 한글 텍스트 스테가노그래피 영역은 향후 연구해야할 부분이다.

## References

- [1] Ji S. S., "A Study of Hangul Text Steganography based on Genetic Algorithm", KIISC, Vol. 21, No. 3, pp. 7-12, 2016.
- [2] Kim J. G. and Jeon J. H., "Factors Affecting Internet User's Information Security Intention: Focused on Computer Virus", KIECA, Vol. 5, No. 2, pp. 47-70, 2011.
- [3] Agarwal M., "Text Steganographic Approaches : A Comparison", International Journal of Network Security&Its Applications, Vol. 5, No. 1, pp. 91-106, 2013.
- [4] Saraswathi V. and Kingslin Sumathy, "Different Approaches to Text Steganography: A Comparison", International Journal of Emerging Research in Management & Technology, Vol. 3, Issue 11, pp. 124-127, 2014.
- [5] Kumar K. A., "A Comparative Result Analysis of Text Based Steganographic Approaches", Journal of Computer Engineering, Vol. 17, Issue 3, pp. 70-74, 2015.
- [6] Nagarhalli T. P., "A New Approach to Text Steganography Using Adjectives", International Journal of Innovative Research in Science, Engineering and Technology, Vol. 4, Issue 5, pp. 3147-3152, 2015.
- [7] Kingslin Sumathy and Kavitha N., "Evaluative Approach towards Text Steganographic Techniques", Indian Journal of Science and Technology, Vol. 8, No. 29, pp. 1-8, 2015.
- [8] Abaas S. T., Z. Khudhair N., Kareem S. K. and Khudhair Zainab Nighal, "Improve Capacity in Text in Text Steganography", European Academic Research, Vol. II, Issue 12, pp. 15049-15061, 2015.
- [9] Yufei Sun, Liangli Ma and Shuang Wang, "A Comparative Evaluation of String Similarity Metrics for Ontology Alignment", Journal of Information & Computational Science, Vol. 12, No. 3, pp. 957-964, 2015.
- [10] Cohen W., Ravikumar P. and Fienberg S. E., "A Comparison of String Distance Metrics for Name-Matching Tasks", American Association for Artificial Intelligence, pp. 73-78, 2003.



지 선 수 (Ji Seon Su)

- 중신회원
- 충남대학교 계산통계학과(학사)
- 중앙대학교 응용통계학과(석사)
- 중앙대학교 응용통계학과(박사)
- 명지대학교 컴퓨터공학과(박사수료)
- (현) 강릉원주대학교 소프트웨어학과 교수
- 관심분야 : 정보보안(암호키, 정보은닉), 스테가노그래피