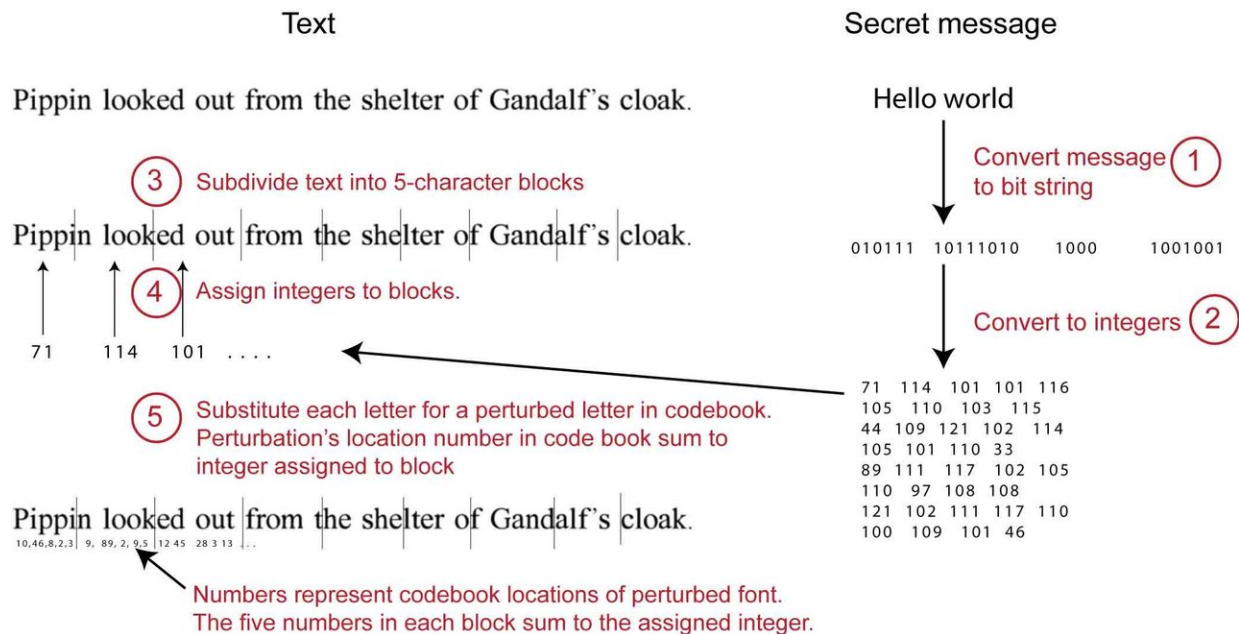


Researchers hide information in plain text

May 10 2018



Someone using FontCode would supply a secret message and a carrier text document. FontCode converts the secret message to a bit string (ASCII or Unicode) and then into a sequence of integers. Each integer is assigned to a five-letter block in the regular text where the numbered locations of each letter sum to the integer. Credit: Changxi Zheng/Columbia Engineering

Computer scientists at Columbia Engineering have invented FontCode, a new way to embed hidden information in ordinary text by imperceptibly changing, or perturbing, the shapes of fonts in text. [FontCode](#) creates font perturbations, using them to encode a message that can later be decoded to recover the message. The method works with most fonts and,

unlike other text and document methods that hide embedded information, works with most document types, even maintaining the hidden information when the document is printed on paper or converted to another file type. The paper will be presented at SIGGRAPH in Vancouver, British Columbia, August 12-16.

"While there are obvious applications for espionage, we think FontCode has even more practical uses for companies wanting to prevent document tampering or protect copyrights, and for retailers and artists wanting to embed QR codes and other metadata without altering the look or layout of a document," says Changxi Zheng, associate professor of computer science and the paper's senior author.

Zheng created FontCode with his students Chang Xiao (Ph.D. student) and Cheng Zhang MS'17 (now a Ph.D. student at UC Irvine) as a [text](#) steganographic method that can embed text, metadata, a URL, or a digital signature into a text document or image, whether it's digitally stored or printed on paper. It works with common font families, such as Times Roman, Helvetica, and Calibri, and is compatible with most word processing programs, including Word and FrameMaker, as well as image-editing and drawing programs, such as Photoshop and Illustrator. Since each letter can be perturbed, the amount of information conveyed secretly is limited only by the length of the regular text. Information is encoded using minute font perturbations—changing the stroke width, adjusting the height of ascenders and descenders, or tightening or loosening the curves in serifs and the bowls of letters like o, p, and b.

"Changing any letter, punctuation mark, or symbol into a slightly different form allows you to change the meaning of the document," says Xiao, the paper's lead author. "This hidden information, though not visible to humans, is machine-readable just as barcodes and QR codes are instantly readable by computers. However, unlike barcodes and QR codes, FontCode doesn't mar the visual aesthetics of the printed

material, and its presence can remain secret."

Data hidden using FontCode can be extremely difficult to detect. Even if an attacker detects font changes between two texts—highly unlikely given the subtlety of the perturbations—it simply isn't practical to scan every file going and coming within a company.

Furthermore, FontCode not only embeds but can also encrypt messages. While the perturbations are stored in a numbered location in a codebook, their locations are not fixed. People wanting to communicate through encrypted documents would agree on a private key that specifies the particular locations, or order, of perturbations in the codebook.

"Encryption is just a backup level of protection in case an attacker can detect the use of font changes to convey [secret information](#)," says Zheng. "It's very difficult to see the changes, so they are really hard to detect—this makes FontCode a very powerful technique to get data past existing defenses."

FontCode is not the first technology to hide a message in text—programs exist to hide messages in PDF and Word files or to resize whitespace to denote a 0 or 1—but, the researchers say, it is the first to be document-independent and to retain the secret information even when a document or an image with text (PNG, JPG) is printed or converted to another file type. This means a FrameMaker or Word file can be converted to PDF, or a JPEG can be converted to PNG, all without losing the secret information.

To use FontCode, you would supply a secret message and a carrier text document. FontCode converts the secret message to a bit string (ASCII or Unicode) and then into a sequence of integers. Each integer is assigned to a five-letter block in the regular text where the numbered codebook locations of each letter sum to the integer.

Recovering hidden messages is the reverse process. From a digital file or from a photograph taken with a smartphone, FontCode matches each perturbed letter to the original perturbation in the codebook to reconstruct the original message.

Matching is done using convolutional neural networks (CNNs). Recognizing vector-drawn fonts (such as those stored as PDFs or created with programs like Illustrator) is straightforward since shape and path definitions are computer-readable. However, it's a different story for PNG, IMG, and other rasterized (or pixel) fonts, where lighting changes, differing camera perspectives, or noise or blurriness may mask a part of the letter and prevent an easy recognition.

While CNNs are trained to take into account such distortions, recognition errors will still occur, and a key challenge for the researchers was ensuring a message could always be recovered in the face of such errors. Redundancy is one obvious way to recover lost information, but it doesn't work well with text since redundant letters and symbols are easy to spot.

Instead, the researchers turned to the 1700-year-old Chinese Remainder Theorem, which identifies an unknown number from its remainder after it has been divided by several different divisors. The theorem has been used to reconstruct missing [information](#) in other domains; in FontCode, researchers use it to recover the original message even when not all letters are correctly recognized.

"Imagine having three unknown variables," says Zheng. "With three linear equations, you should be able to solve for all three. If you increase the number of equations from three to five, you can solve the three unknowns as long as you know any three out of the five equations."

Using the Chinese Remainder theory, the researchers demonstrated they

could recover messages even when 25% of the letter perturbations were not recognized. Theoretically the error rate could go higher than 25%.

The authors, who have filed a patent with Columbia Technology Ventures, plan to extend FontCode to other languages and character sets, including Chinese.

"We are excited about the broad array of applications for FontCode," says Zheng, "from [document](#) management software, to invisible QR codes, to protection of legal documents. FontCode could be a game changer."

The study is titled "FontCode: Embedding Information in Text Documents using Glyph Perturbation."

More information: Chang Xiao et al, FontCode, *ACM Transactions on Graphics* (2018). [DOI: 10.1145/3152823](https://doi.org/10.1145/3152823) , arxiv.org/pdf/1707.09418.pdf

Provided by Columbia University School of Engineering and Applied Science

Citation: Researchers hide information in plain text (2018, May 10) retrieved 20 March 2024 from <https://techxplore.com/news/2018-05-plain-text.html>

This document is subject to copyright. Apart from any fair dealing for the purpose of private study or research, no part may be reproduced without the written permission. The content is provided for information purposes only.
