

Development of an OCR System

Nathan Harmata

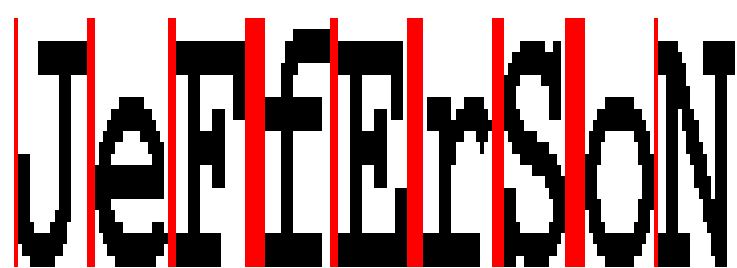
TJHSST Computer Systems Lab 2007 - 2008

Abstract

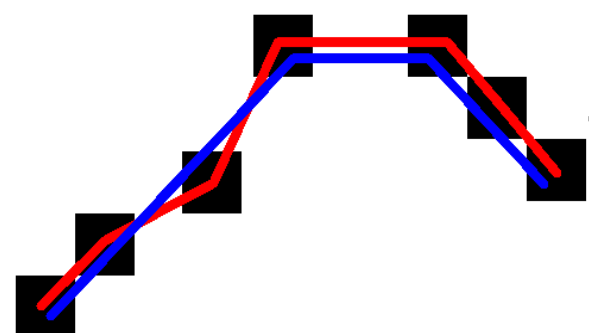
OCR (Optical Character Recognition) is a very practical field of Computer Science. Since the late 1980's, researchers have been developing system to identify text from non electronic sources, such as pictures or newspapers. The use of OCR systems has spanned from making books in Braille to sorting mail by zip code.

Procedures

The input for the current prototype is a PNG picture file that contains text in the Courier font. Using the Java BufferedImage class, locations and colors of the pixels in the image can be determined. The program uses these to find the positions of horizontal straight lines of whitespace in the image. It pairs together lines of whitespace and ignores those so that only lines of text remain. Each line is parsed into words using a similar method involving vertical lines. After making spacing analysis, each word is parsed into letters, as seen below.

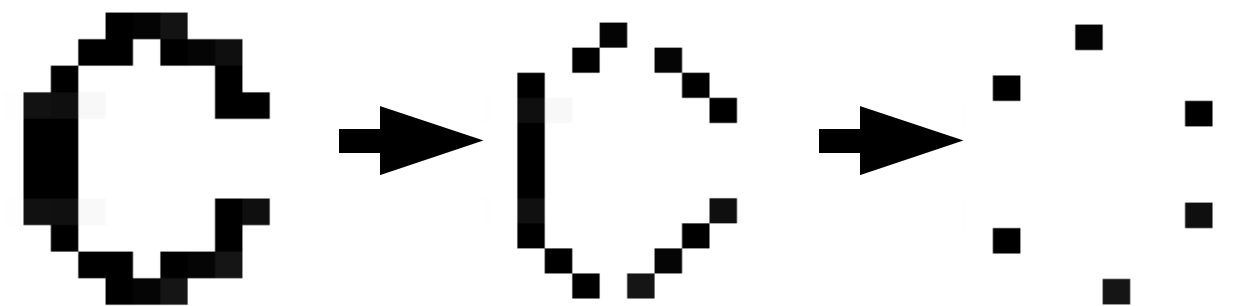


Then, each letter is parsed into portions that pass the vertical line test. After getting rid of unnecessary information, each portion is then transformed into a sum of line segments of different slopes.



Background

Although there are a few options currently available to the public, like Microsoft Document Imaging, most of them are either unused or not accurate enough. The goal of this project is to create an OCR system that is simple to use and can handle most formatting and fonts.



The result is that a letter is simplified into a few line segments connecting key pixels. This procedure is applied to each letter of several different fonts, and information from the results is stored and averaged. Using this, a cache is created to which results from OCR analysis can be dynamically compared.

Results

The groupings of letters that share the same cache value are shown to the right. Using these and a dictionary, a list of all possible matching words can be generated. As of now, the results are not spread out enough to yield consistent word recognition.

a
eo
g
cdq
sz
lnr
it
fj
hmu
p
v
xy
bk
w

The goal of 3rd Quarter will be to develop more heuristics with which to differentiate letters in the cache so that the results are more spread out.