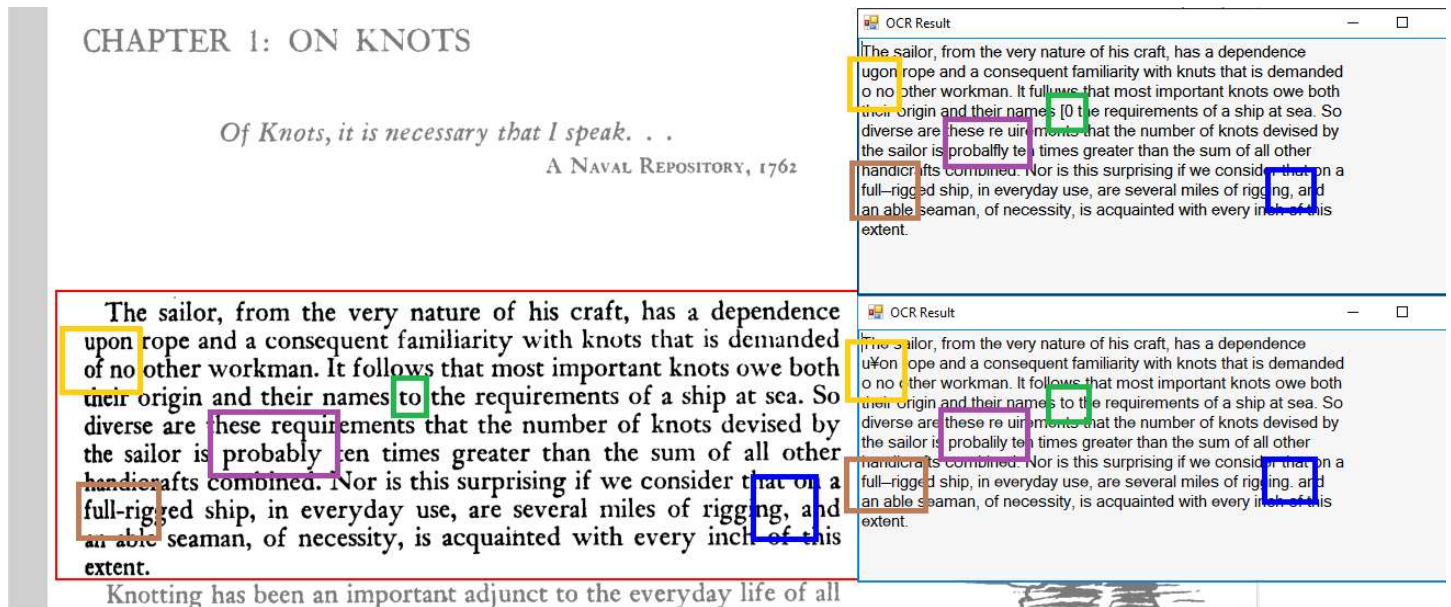


OCR & R2V

Undoubtedly many users will come across the need to run OCR (Optical Character Recognition) or need a R2V (Raster to Vector) convertor in order to either make a scanned document searchable or a drawing editable.

Many OCR solutions tend to use the Google Tesseract Engine

It is not 100% but in this "Snipping-Ocr" sample of a screen snip of 627 characters (including white space) it made only a few understandable errors, highlighted here in coloured boxes (red is the original on screen)



The diagram shows the source on left and two separate attempts to convert screen to text and here is the result

The first line seems perfect

The sailor, from the very nature of his craft, has a dependence

The next two lines have characters in the second position on the line that overlap (see the orange box above) this failure is understandable and may be avoided if we zoom in.

u#on rope and a consequent familiarity with knots that is demanded
o no other workman. It follows that most important knots owe both

The following line was wrong on first run (see top right green box) by zooming in was correct on second run.

their origin and their names to the requirements of a ship at sea. So

The fifth and sixth lines (same as the second and third) also suffered where two characters q and l overlapped

diverse are these re uirements that the number of knots devised by
the sailor is probalily ten times greater than the sum of all other

7th was good

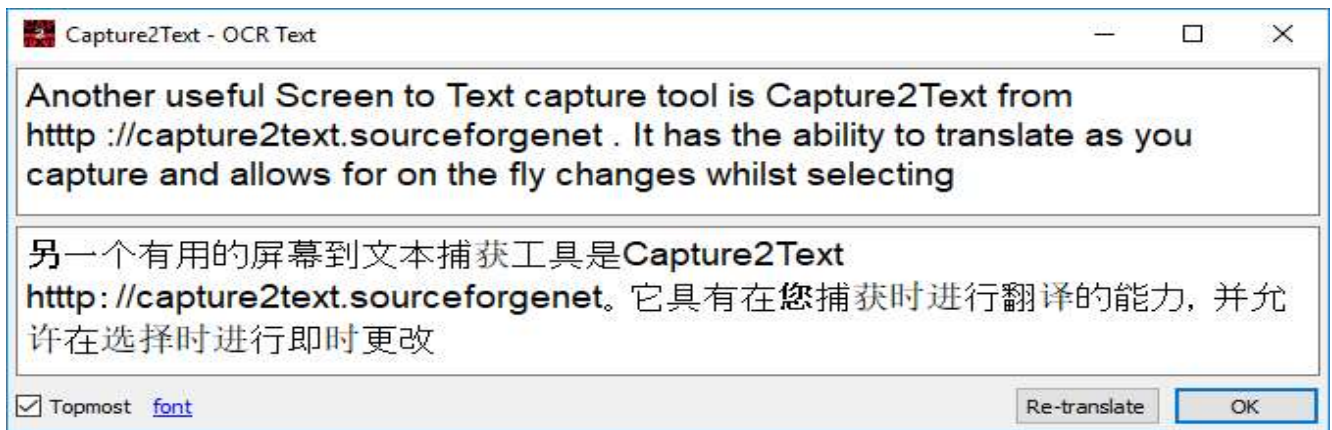
handicrafts combined. Nor is this surprising if we consider that on a

The 8th line has 2 different punctuation issues the hyph-en was converted to an em and in the 2nd run , is .

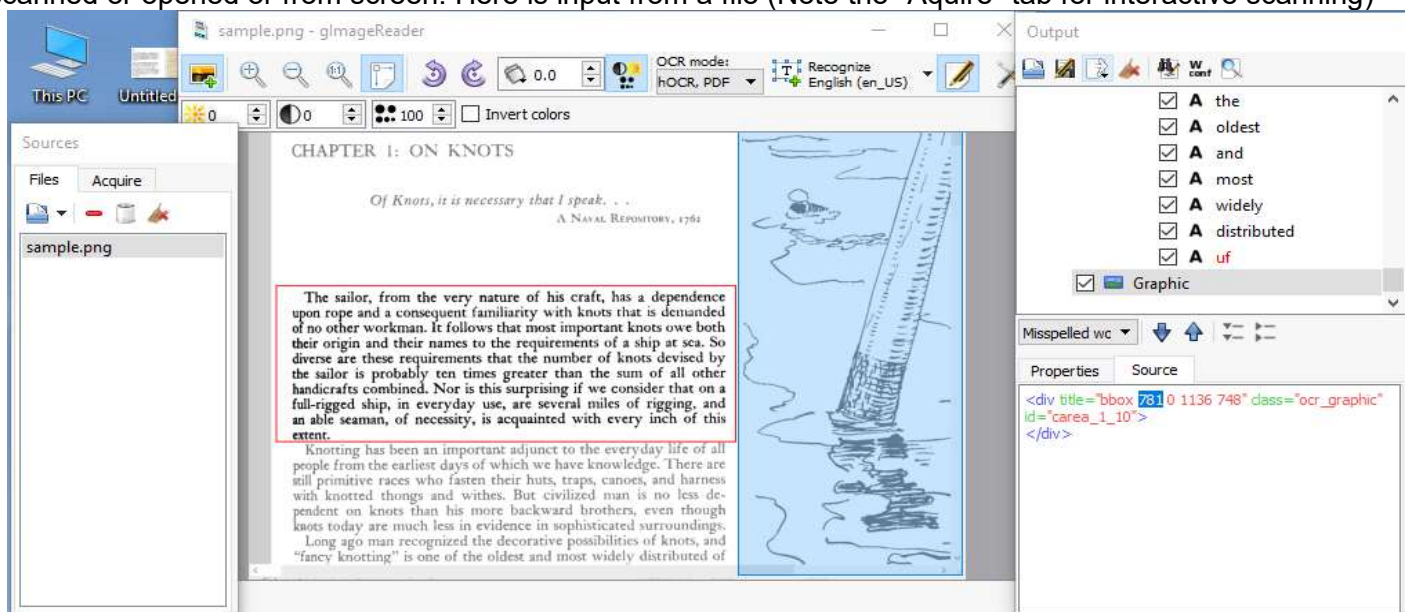
full—rigged ship, in everyday use, are several miles of rigging. and
an able seaman, of necessity, is acquainted with every inch of this extent.

I consider those minor issues that can be easily corrected and at 99.2% success rate (5 wrong out of 627) is good the source and binary is at <https://github.com/thepirat000/Snipping-Ocr/releases>

Another useful Screen to Text capture tool is Capture2Text from <http://capture2text.sourceforge.net/> It is able to translate as you capture and allows on the fly changes such as the full stop in softforge.net



However a favourite for text and graphics is <https://github.com/manisandro/gImageReader>. Output can be Html ,PDF or Text and it has floating pallets to maximise screen use that allow interactive editing of files scanned or opened or from screen. Here is input from a file (Note the "Aquire" tab for interactive scanning)



Resultant PDF with OCR'd text and a graphic for good measure

