

Kraken - an Universal Text Recognizer for the Humanities

Benjamin Kiessling (benjamin.kiessling@psl.eu), Université PSL, France; Leipzig University

1. Introduction

Retrodigitization of both printed and handwritten material is a common prerequisite for a diverse range of research questions in the humanities. While optical character recognition on printed texts is widely considered to be fundamentally solved in academia, with the most commonly used paradigm (Graves et al., 2006) dating back to 2006, this hasn't translated into increased availability of adaptable, libre-licensed OCR engines to the technically inclined humanities scholar.

The nature of the material of interest commands a platform that can be altered with minimum effort to achieve optimal recognition accuracy; uncommon scripts, historical languages, complex or archaic page layout, and non-paper writing surfaces are rarely satisfactorily addressed by off-the-shelf commercial solutions. In addition, an open system ameliorates the severe resource constraints of humanities research by enabling sharing of artifacts, such as training data and recognition models, inaccessible with proprietary OCR technology.

2. Kraken

The Kraken text recognition engine is an extensively rewritten fork of the OCRopus system. It can be used both for handwriting and printed text recognition, is easily (re-)trainable, and great care has been taken to eliminate implicit assumptions on content and layout that complicate the processing of non-Latin and non-modern works.

Thus Kraken has been extended with features and interfaces enabling the processing of most scripts, among them full Unicode right-to-left, bidirectional, and vertical writing support, script detection, and multiscript recognition. Processing of scripts not included in Unicode is also possible through a simple JSON interface to the codec mapping numerical model outputs to characters. The same interface provides facilities for efficient recognition of large logographic scripts.

Output includes fine-grained bounding boxes down to the character level that may be used to quickly acquire a large number of samples from a corpus to assist in paleographic research. Kraken implements a flexible output serialization scheme utilizing a simple templating language. Templates are available for the most commonly used formats ALTO, hOCR, TEI, and abbyyXML.

While including implementations of all the subprocesses needed in a text recognition pipeline, most functional blocks can be accessed separately on the command line, allowing flexible substitution of specially optimized methods. A stable programming interface allows total customization and integration into other software packages.

3. Recognition

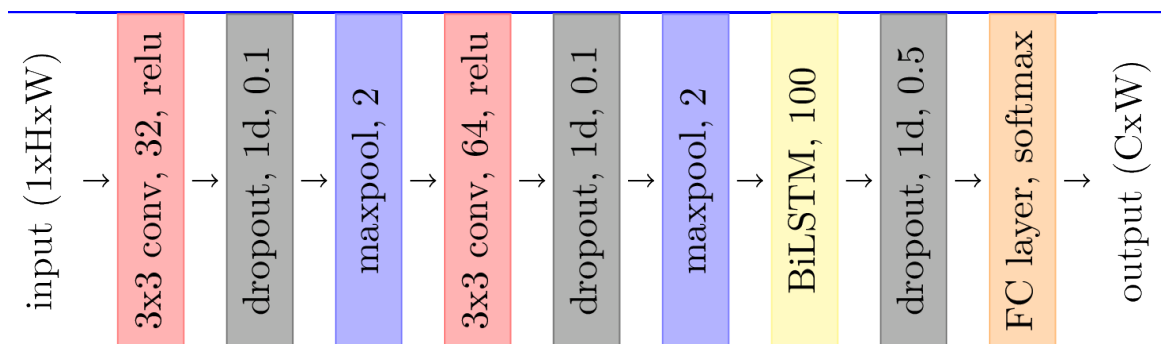


Figure 1. Network architecture (H : sequence height, W : sequence length, C : alphabet size)

The recognition engine operates as a segmentation-less sequence classifier using an artificial neural network to map an image of a single line of text, the input sequence, into a sequence of characters, the

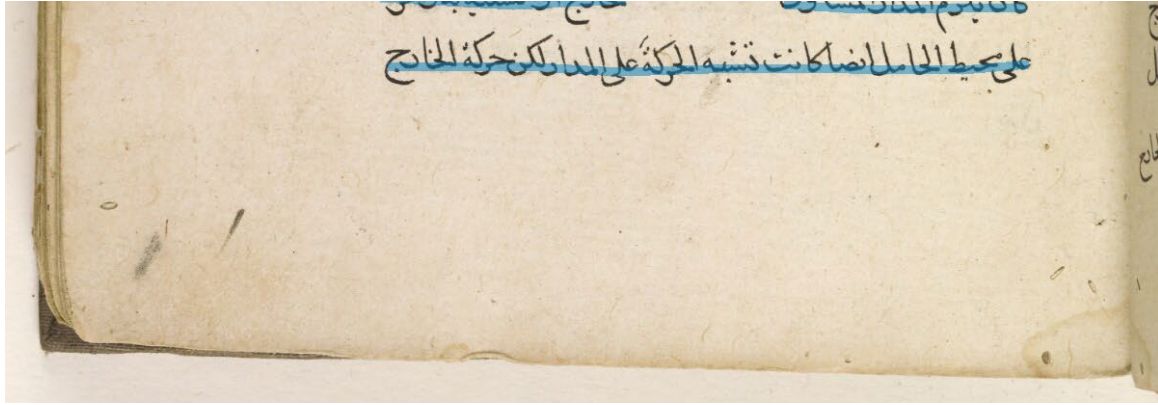


Figure 2. Sample output of the trainable segmentation method.

Kraken's layout analysis extracts text lines from an input image for later processing by the recognition engine. Apart from a basic segmenter taken from OCRopus a trainable line extractor is in the process of being implemented. Full trainability of layout analysis is of utmost importance to a truly universal OCR system, as text layout and its semantics varies widely across time and space, e.g. hand-crafted methods for printed Latin text are unlikely to work reliably on Arabic text or manuscripts with extensive interlinear annotation.

The trainable layout analysis module consists of a two-step instance segmentation method: an initial seed-labelling network operates on the whole page labelling the area between baseline and mean of each line. As the output of the network is a probability of each pixel belonging to a baseline it is binarized using hysteresis thresholding after smoothing with a gaussian filter. The binarized image is then skeletonized and end point are extracted with a discrete convolution. Finally, the vectorized baseline between the endpoints is rectified and a variable environment calculated based on the distance of connected components from the labelled area is extracted.

The seed-labelling network is a modified U-net (Ronneberger et al., 2015) on the basis of a 34-layer residual network (He et al., 2016) pretrained on ImageNet.

Preliminary results on a page from a publicly available dataset of Arabic and Persian manuscripts (Kiessling et al., 2019) can be seen in Figure 2.

Script detection, the basis for multi-script support in the recognizer, is implemented as a segmentation-less sequence classification problem, similar to text recognition. Instead of assigning a unique label to each code point or grapheme cluster we assign all code points of a particular script the same label. The network is then trained to output the correct sequence of script labels (Figure 3). The output sequence is then used to split the line into single-script runs that can be classified with monolingual recognition models (Figure 4).

ويقول رئيس شركة U. S. Steel : « انا لا أومن بالدين . فالدين له الميزة الغير المسرة

ويقول رئيس شركة U. S. Steel : « انا لا أومن بالدين . فالدين له الميزة الغير المسرة

00000020000020000002002000000222000000200002002000222212212211112000020000200000

Figure 3. Original and modified ground truth (top: original line, middle: transcription, bottom: assigned script classes)

5. Results

	Mean character accuracy	Standard deviation	Maximum accuracy
Prints			
Arabic (Kiessling et al., 2017)	99.5%	0.05	99.6%
Persian ¹	98.3%	0.33	98.7%
Syriac ²	98.7%	0.38	99.2%
Polytonic Greek ³	99.2%	0.26	99.6%

Latin (Springmann et al., 2018)	98.8%	0.09	99.3%
Latin incunabula (Springmann et al., 2018)	99.0%	0.11	99.2%
Fraktur (Springmann et al., 2018)	99.0%	0.31	99.3%
Cyrillic Manuscripts	99.3%	0.15	99.6%
Hebrew 4	96.9%	-	-
Medieval Latin 5	98.2%	-	-

اجواء شتى . وهذا سبب من اسباب الغموض اللغوي وعدم تضمين اللفظ معناه المحدود الوضعي ، والداعي الى تفوق الفن على الطبيعة ليصبح ضرباً من ضروب المخدرات . وكذا فلا يُصور الشيء بكامله وانما يُكتفى بالاشارة الى بعض اجزائه .

وخاصة الاليجاء قائمة على امكانياته ، لان الاليجاء كالعنسة يوسع الافتراضات ويطلق الخيلة والشعور والتأمل . فالاليجاء يحوِّك ألوفاً من الخيوط الملونة حول مغزل النفس ، او قل انه انتظار شيء سيحدث ؛ وفي الانتظار لذة لا تعرفها الحقيقة الواقعة . وكأما تنقسم الالفاظ الى نثرية وشعرية . والشعرية هي الغرض . ومن هذا القبيل فن التصوير الزيتي ، حيث يولد انسجام الالوان ، وازلالها ، وادهان بعض الخطوط والاجزاء ، في نفس المتشبع ، جواً ايجائياً شبيهاً بتأثير الالفاظ الشعرية . وفن « رمبراندت » هو من هذا الباب ، حيث تحمل اللوحة اليك ما ينطلق وراء الخطوط والالوان ، وفي هذا الاتجاه يقول احدهم :

«Deux choses sont également requises : l'une est certaine somme de complexité ou plus proprement de combinaison; l'autre une certaine quantité d'esprit suggestif, quelque chose comme un courant souterrain de pensée non visible, indéfinie... c'est l'excès dans l'expansion du sens qui ne sait être qu'insinué. » (١)

وللمشترع ملارمه قول صراح بهذا الصدد : « ان ما شيد من صروح ، والبحر ، والوجه الانساني ، متعة لا يؤديها الوصف ، بل يجملها الاليجاء »^(٢) . واذن ففي الاليجاء غنى للفكر « لان فعل اللغة الاليجائية الرئيسي هو توليد مجاري فكرية وشعرية ، تقاربت ام تفرعت . »^(٣)

والشعر الرمزي بتحديد شحنة ايجائية . « ففي هذا اليجاء يتراكم السكوت ويحفظ ويمتد . هذا ما يبدو في ابتسامة « الجوكوندا » وفي « اصبع يوحنا » المرفوعة في تمثال دي فنشي . ثم ان الاليجاء بمثابة عصا الجوقة الموسيقية . فكما ان سيلاً من الموسيقى يتدفق تحت شارته ، كذا تنبجس الدنيا الداخلية بانفعال اللفظة الاليجائية .

(١) Bremond, La Poésie pure, p. 118.

(٢) Mallarmé, Divagations, p. 245.

(٣) Paulhan, La Double Fontc. du lang. P. 169

Kraken has been used on a wide variety of writing systems, achieving uniformly high character accuracy (CER). Sample accuracies for a diverse set of scripts spanning across multiple centuries of printing are shown in Table 1.

As a special use case we evaluated recognition of text and emphasis in a mixed English and romanized Arabic library catalog on a training set of 350 lines (50 lines in the validation set) resulting in an averaged CER of 99.3% ($\sigma=0.16$) over 10 runs with 95.38% CER on cursive and text with increased spacing ($\sigma=1.46$). When using only emphasized text accuracy as the stopping criterium mean accuracy rises to 99.03% ($\sigma=0.28$).

Appendix A

Bibliography

1. **Graves, A., Fernández, S., Gomez, F. and Schmidhuber, J.** (2006). Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. *Proceedings of the 23rd International Conference on Machine Learning* . ACM, pp. 369–376.
2. **He, K., Zhang, X., Ren, S. and Sun, J.** (2016). Deep residual learning for image recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* . pp. 770–778.
3. **Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I. and Salakhutdinov, R. R.** (2012). Improving neural networks by preventing co-adaptation of feature detectors. *ArXiv Preprint ArXiv:1207.0580* .
4. **Kiessling, B., Miller, M. T., Maxim, G., Savant, S. B. and others** (2017). Important New Developments in Arabographic Optical Character Recognition (OCR). *Al-'Uṣūr Al-Wuṣṭā* , **25** : 1–13.
5. **Kiessling, B., Stoekl Ben Ezra, Daniel and Miller, Matthew Thomas** (2019). BADAM: A Public Dataset for Baseline Detection in Arabic-script Manuscripts.
6. **Ronneberger, O., Fischer, P. and Brox, T.** (2015). U-net: Convolutional networks for biomedical image segmentation. *International Conference on Medical Image Computing and Computer-Assisted Intervention* . Springer, pp. 234–241.
7. **Springmann, U., Reul, C., Dipper, S. and Baiter, J.** (2018). Ground Truth for training OCR engines on historical documents in German Fraktur and Early Modern Latin. *ArXiv Preprint ArXiv:1809.05501* .

Notes

1. Mid-20th century printing
 2. Late-19th century printing in Serṭā form
 3. Late-19th century printing
 4. Midrash Tanhuma, BNF Héḇ 150
 5. Josephus Latinus, Bamberg 78 with augmentation
-