# Call for Bachelor Thesis

**FIZ** Karlsruhe
Leibniz Institute for Information Infrastructure

## *Handwritten and Printed Text Separation in Historical Documents*

### Objective of this work:

With the increase of digitized documents, automatic document analysis has become extremely important. The presentation of historical documents to the public introduces a variety of document types, content, quality and structure. Fundamentally speaking, documents can be skewed, noisy, and overlapped with graphics, i.e., lines, unconstrained annotations, stamps.
Most optical character recognition (OCR) systems recognize either printed or handwritten text. Hence, the task of the thesis is to separate machine printed text from handwritten text in scanned documents before feeding it to an OCR system.

In this thesis:

1. Documents containing a mix of handwritten and printed text will be collected.
2. An additional mixed dataset may be generated from historical documents.
3. The existing approaches of text separation will be reviewed and investigated.
4. A pixel-based approach for text separation based on [1] will be applied.
5. The results will be evaluated based on the ground truth data.

[1] Dutly, N., Slimane, F., & Ingold, R. (2019, September). Phti-ws: A printed and handwritten text identification web service based on fcn and crf post-processing. In *2019 International Conference on Document Analysis and Recognition Workshops (ICDARW)* (Vol. 2, pp. 20-25). IEEE.
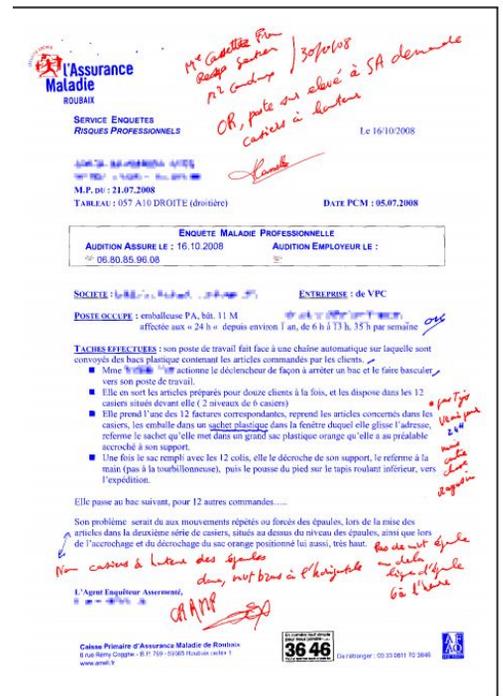
The project work will be supervised by **Prof. Dr. Harald Sack, Mahsa Vafaie and Oleksandra Bruns**, **Information Service Engineering at Institute AIFB, KIT, in collaboration with FIZ Karlsruhe.**

**Keywords:**
Machine Learning, CNN, pattern recognition
**Pre-requisites:**
Knowledge of Programming with Python.

Contact persons:
**Mahsa Vafaie**
mahsa.vafaie@kit.edu

**Oleksandra Bruns**
oleksandra.bruns@kit.edu

**Institute of
Applied Informatics and
Formal Description Methods**
http://www.aifb.kit.edu/

**KIT**
Karlsruhe Institute of Technology