

Ancient Document Text Recognition: LLM OCR Agent

Seorin Kim & Julien Baudru

Abstract

Your abstract.

1 Introduction

The aim of this pilot project is to see whether generative AI methods, and multimodal LLMs in particular, can be used effectively to detect text in old administrative documents. All the documents studied here are produced in handwriting that is difficult for humans to decipher, and the manual transcription process tends to take a considerable amount of time. We compare traditional OCR/HTR text detection methods with the results obtained by an LLM and those obtained by manual annotation.

2 Similar works

In this section, we present a summary of the previous research aiming to combine text detection with language models (Gpt, Bert, etc.).

In [Fuj23], the authors introduce a new method for text recognition called Decoder-only Transformer for Optical Character Recognition (DTrOCR). DTrOCR uses only a decoder, using a pre-trained generative language model, in contrast to traditional encoder-decoder methods. The authors tested whether a successful natural language processing model could be applied to text recognition in computer vision. Their experiments showed that DTrOCR significantly outperformed current state-of-the-art methods in recognising printed, handwritten and scene text in both English and Chinese.

In [TGL24], the authors address the challenge of poor OCR quality in digitised historical documents, which is a barrier to humanities research. Traditional post-OCR correction methods use sequence-to-sequence models. Instead, the authors propose the use of genera-

tive language models with a prompt-based approach. By tuning Llama 2 with prompts and comparing it to a fine-tuned BART model on 19th century British newspaper articles, they demonstrate significant improvements in OCR error correction. Llama 2 achieves a 54.51% reduction in the character error rate, outperforming BART's 23.30% reduction. This approach shows promise for improving the accessibility of historical texts for researchers.

In [Lö23], the author introduces a method to digitize over 100,000 historic plans from the Swiss Archive for Landscaping Architecture using AI models. The approach employs a three-model architecture: a layout model to identify text, an OCR model to extract words, and a named entity recognition (NER) model to label key information. K-means clustering groups text blocks for OCR processing. Various deep-learning models were evaluated, including German BERT for NER, and retrained on the NVIDIA DGX-2 system. The pipeline achieved an F1 score of 48%, with the NER model scoring 86% and the OCR model correctly extracting 54% of words.

LLM here are multimodal, so for the image the model use DNN to decrypt image, cite some existing model using it

3 Framework

4 Experiments

- Compare LLM vs OCR
- Compare LLM vs HTR
- Compare LLM vs Human
- Ask LLM to correct extracted text

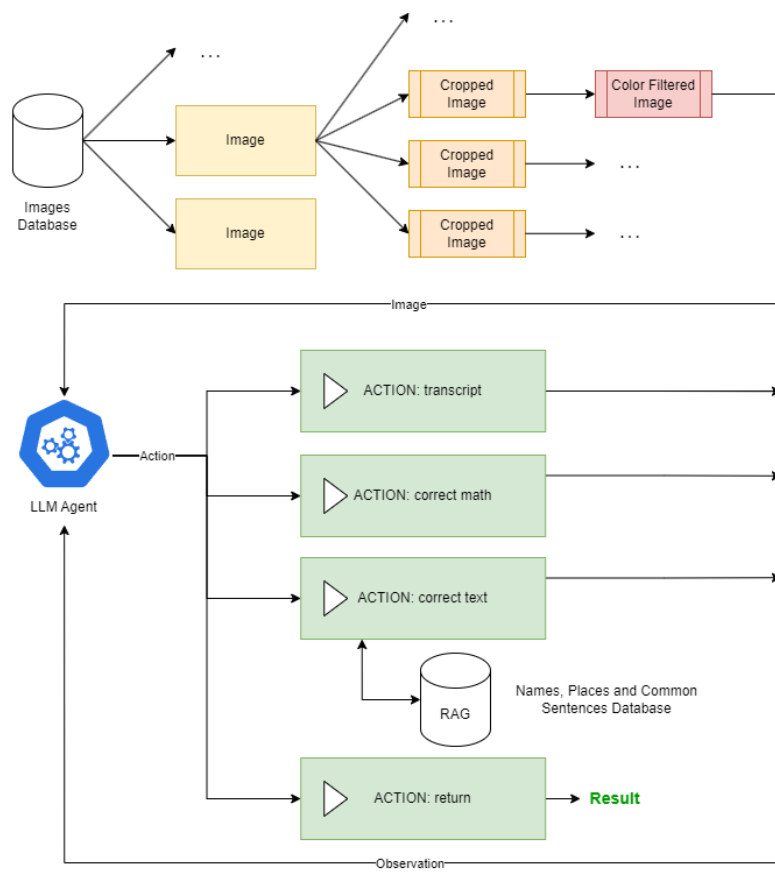


Figure 1: LLM OCR agent

References

- [Fuj23] Masato Fujitake. Dtrocr: Decoder-only transformer for optical character recognition, 2023.
- [Lö23] Kevin Löffler. *Digitize Historic Architectural Plans with OCR and NER Transformer Models*. Other thesis, OST Ostschweizer Fachhochschule, May 2023. Thesis advisor: Mitra Purandare.
- [TGL24] Alan Thomas, Robert Gaizauskas, and Haiping Lu. Leveraging LLMs for post-OCR correction of historical newspapers. In Rachele Sprugnoli and Marco Passarotti, editors, *Proceedings of the Third Workshop on Language Technologies for Historical and Ancient Languages (LT4HALA) @ LREC-COLING-2024*, pages 116–121, Torino, Italia, May 2024. ELRA and ICCL.