

CDD Ing nieur R&D - Computer Vision Machine Learning Tours - LIFAT

Contexte

Ce poste s'inscrit dans le cadre du projet ANR TypoRef dirig  par Remi Jimenes (CESR Tours). Con u dans une approche interdisciplinaire, le projet vise   inventorier, d crire et  tudier les mat riaux typographiques fran ais de la Renaissance en d veloppant une plateforme et des outils d'analyse et d'exploration de corpus de livres num ris s.

La num risation des livres anciens a d but  au milieu des ann es 1990. Depuis cette date, des  quipes de chercheurs alliant informaticiens et sp cialistes du patrimoine  crit ont collabor  pour d velopper des outils d di s. Longtemps, l'analyse et la reconnaissance d'images de documents ont  t  consid r es comme des t ches complexes d compos es en 2 sous-t ches ind pendantes : l'analyse de la mise en page et la reconnaissance du texte. L'analyse de la mise en page vise   d tecter tous les composants de l'image tels que les blocs de texte, les tableaux, les images, les graphiques, ... La reconnaissance de texte est consacr e   la reconnaissance des caract res   l'int rieur des blocs de texte.

Dans le cadre d'une collaboration  tablie en 2003, le CESR et le LIFAT ont d velopp  les logiciels [Agora and Retro](#). Les d veloppements ont  t  prolong s en 2011 et 2012 par deux bourses successives "Award in digital Humanities " financ es par Google. Ces outils d'analyse et d'indexation de la mise en page adapt s aux images des livres imprim s de la Renaissance ont permis d'extraire et d'indexer automatiquement quelque 13 500 gravures, lettrines, bandeaux et fleurons   partir de 600 fac-simil s des biblioth ques virtuelles humanistes du BVH.

Agora et Retro ont d j  plus de 15 ans et sont bas s sur des technologies aujourd'hui d pass es au vu des d veloppements de ces derni res ann es. R cemment, des architectures profondes ont am lior  les performances obtenues pour les deux sous-t ches. Sur la base de la tr s forte exp rience accumul e au cours des vingt derni res ann es, le projet TypoReF entend permettre le d veloppement de nouveaux outils d'analyse de la mise en page, d'identification et d'indexation des mat riaux typographiques adapt s aux livres imprim s anciens, en utilisant les technologies les plus puissantes du moment, notamment celles bas es sur l'apprentissage profond.

Missions   r aliser

L'objectif est la mise en place d'une plateforme web de labellisation de contenus. Cette labellisation concerne   la fois la collecte de m tadonn es produites par un utilisateur expert et la production de m tadonn es r sultant d'une analyse des images/pixels. Cette double indexation (humain/machine) doit permettre de rechercher, comparer, regrouper des formes, et ainsi de mettre en  vidence des liens difficilement perceptibles.

La plateforme n cessite de d velopper les fonctionnalit s suivantes :

- Segmentation s mantique des images de documents   l'aide d'architectures d'apprentissage profond (pour remplacer l'ancienne Agora) : localisation et caract risation d' l ments de contenu sp cifiques   l'int rieur des pages num ris es d'un livre imprim  ancien (blocs de texte, d cor grav ...) avec diff rents niveaux de d tail (au sein d'un bloc de texte, on distinguera ainsi des lignes, des caract res) et caract risation de ces  l ments (lettrage, bandeau, marque, caract res, etc.) Chaque  l ment sera associ    un ensemble de m tadonn es permettant de le d crire et de retrouver pr cis ment son origine.
- Le clustering des  l ments de contenu extraits (pour remplacer l'ancien R tro) consiste   faire des comparaisons entre les  l ments, principalement des ornements grav s (comme les lettrages ou les bandeaux), pour proposer des correspondances entre des formes similaires. La plateforme utilisera diff rents types d'algorithmes d'apprentissage automatique, notamment l'apprentissage dit " non supervis  ".
- Importation/exportation des donn es depuis/vers la base de donn es selon des standards tels que ALTO ou IIIF .

Comp tences

- D veloppement logiciel et ing nierie, Python, HTML, CSS, JS
- Apprentissage automatique, vision par ordinateur
- Capacit    travailler en  quipe, esprit curieux et rigoureux.

Comment postuler ?

- Envoyez votre CV et votre lettre de motivation
- Contact : Thierry Brouard, Jean-Yves Ramel --> prenom.nom AT univ-tours.fr
- Lieu : LIFAT et CESR - Tours France - <http://lifat.univ-tours.fr>
- Poste : (environ) contrat de 12 mois - dur e et salaire selon l'exp rience. Temps plein
- Dates du contrat : 1er septembre 2023 - 31 ao t 2024

Engineer R&D - Computer Vision - Machine Learning Tours - LIFAT

Context

This position is part of the ANR TypoRef project directed by Remi Jimenes (CESR Tours). Conceived in an interdisciplinary approach, the project aims at inventorying, describing and studying French typographic materials of the Renaissance by developing tools for the analysis and exploration of digitized book corpora and by setting up a dedicated digital platform.

The digitization of old books began in the mid-1990s. Since then, teams of researchers combining computer scientists and specialists in written heritage have collaborated to develop tools dedicated to the exploitation of these digital corpora. For many years, document image analysis and recognition were considered complex tasks decomposed into two independent sub-tasks: layout analysis and text recognition. Layout analysis aims at detecting all image component such as text blocks, tables, images, graphics, signatures... Text recognition is devoted to the recognition of characters inside text blocks (paragraphs).

In the framework of a collaboration established in 2003, the CESR and the LIFAT have developed the [Agora and Retro](#) software. The developments were extended in 2011 and 2012 by two successive "Digital Humanities" grants funded by Google. These layout analysis and indexing tools adapted to the images of printed books of the Renaissance have made it possible to automatically extract and index some 13,500 engravings, lettering, headbands and fleurons from 600 facsimiles of the BVH Humanist Virtual Libraries.

Agora and Retro are already more than 15 years old and based on technologies that are now outdated given the developments of recent years. Recently, deep architectures have improved the state-of-the-art performance for both sub-tasks. Based on the very strong experience accumulated during the last twenty years, the TypoReF project intends to allow the development of new tools for the analysis of layout, identification and indexing of typographic materials adapted to the ancient, printed book, using the most powerful technologies of the moment, in particular those based on deep learning.

Missions to achieve

The goal is **the implementation of a web-based content labeling platform**. This labeling concerns both the collection of metadata produced by an expert user and the production of metadata resulting from an analysis of the pixels, their layout, and the relationships between the extracted elements of content. This double indexing (human/machine) must allow to search, compare, and group shapes, and thus to highlight links that would otherwise be difficult to perceive.

The platform requires developing the following functionalities:

- **Semantic segmentation of document images using deep learning architectures** (to replace old Agora). It consists of the localization and characterization of specific elements of content inside the digitized pages of an old printed book (text blocks, engraved decoration, etc.) with different levels of detail (within a text block, we will thus distinguish lines and characters). Characterization of these elements is also needed: lettering, bands, marks, characters, etc. A set of associated metadata allows one to characterize and retrieve precisely the location of each component.
- **Clustering of the extracted elements of content** (to replace old Retro) consists in making comparisons between elements, mainly engraved ornaments (such as lettering or headbands), to propose matches between similar forms. The platform will use different types of **machine learning algorithms, including the so-called "unsupervised" learning**.
- **import/export data from/to the internal database according to standards such as ALTO or IIIF** .

Skills

- Software programming and software engineering in Python and web (HTML, CSS, and JS)
- Machine Learning, Computer vision
- Ability to work in a team, curious and rigorous spirit

How to apply?

Send CV and motivation letter

- **Contact:** Thierry Brouard, Jean-Yves Ramel --> prenom.nom AT univ-tours.fr
- **Location:** LIFAT and CESR –Tours France - <http://lifat.univ-tours.fr>
- **Position:** (around) 12-month contract - duration and salary according to experience.
- **Dates of the contract:** Sept 1st 2023 – August 31st 2024