

# Layout segmentation and classification of visual style elements in newspapers

- **Topic:** data science, data mining, machine learning
- **City and country:** Sophia-Antipolis, France
- **Team or project in the lab:** Centre de Recherche Inria Sophia Antipolis – Méditerranée, Biovision Lab
- **Name and mail of the advisors:** Hui-Yin Wu (hui-yin.wu@inria.fr), Pierre Kornprobst (pierre.kornprobst@inria.fr)

## General presentation of the topic

For the adaptation of print to digital journalism, a necessary step is the identification of the various visual elements in a newspaper page, with varied applications such as curation and accessibility. This is an extremely complex challenge due mainly to (1) the the large number visual (e.g., images, logos, ads) and textual (e.g., titles, headings, captions, cross-references), (2) the multi-heirarchical organization of elements (e.g., an article contains headings, images, columns, which contain paragraphs etc.), and (3) the different style constraints from one newspaper to another [2]. To date, existing approaches are only able to identify a small subset of these elements, such as differentiating text from images, and do not address at all the question of hierarchy [1,3,4]. A robust approach to segmenting and classifying newspaper elements thus requires incorporating knowledge on visual design and style into the algorithms for this task.

## Objective of the TER

The objectives of this TER is to improve a current workflow for segmenting and classifying text and visual elements of newspapers using the understanding of visual styles and layout. To achieve this goal, this project will involve :

1. familiarize with the existing annotation tools and propose solutions to automate and facilitate the creation of datasets,
2. investigate machine learning approaches to instance segmentation and classification, notably Mask-RCNN, for complex document segmentation, and train and iteratively improve the model, and
3. establish a baseline on which we can compare our approach to state-of-the-art results.

The work will be carried out on the basis of a pre-existing framework that includes annotation, segmentation, and classification [5]. Training and testing will be conducted on a set of segemented and annotated newspapers.

