# Proposal for DSPRO2 (FS24) – Reading French

## Group Members

Simon Immer,
Nicola Le,
Lars Kehrle

## Short Project Description

We want to develop a model that recognizes diacritic symbols in words. For this project we will focus on the acute accent (◌́), circumflex accent (◌̂) and grave accent (◌̀).
Internal analyses at Post finance have shown that current optical character recognition (OCR) models struggle with detecting diacritics.
Therefore, our model would enrich certain characters with the correct diacritic to improve the accuracy of the recognition result.

## Data Description

Because of a lack of production data and privacy concerns we will create a generator. The generator will create various images containing words. The images will be artificially noisy to simulate scanning artifacts.
The generator will also generate a region of interest that contains the diacritic.
Because we are generating the data in this manner, we will not need to label the data.

## Cloud Service Integration

We will be training our model on our local computers because we do not have the resources to execute everything on the cloud.
We have decided to use Azure devops, in particular the Azure Artifacts service, to publish our trained models.

In industry it is common to publish each model after the training to easily access and deploy.

The free pricing tier should allow us to fulfill our needs with 2GB of storage.
Azure DevOps Services Pricing | Microsoft Azure

## Kanban Tool

For this project, we will use GitHub Projects to track our progress and the remaining work.

# Experiment Tracking Tool Approach

As our experiment tracking tool, we will use Weights and bias DB.

In the experiment tracker, we log the model and the performance, after which we can decide which model best fulfills the requirements.
After that we can easily download our model because with each run, we will keep track of the azure artifact id.