# *** - Registration certificate readout

**Background on the problem**
- Customers upload registration certificates as a scan into ***
- Currently, these documents are read out with a lot of manual work
- They are used as a part of the vehicle procurement process and a faster read-out leads to a faster wow-moment of the customer because the vehicle is ready earlier and there is no suspense.
- We aim to read out three values from german registration certificates:
  - Registration date (3)
  - VIN (Vehicle identification number) (2)
  - License plate (1)



**Goal**

The goal is to evaluate if registration certificates can be read out automatically with good confidence: It is ok if not all the files can be processed fully automatically for all fields, but it is important that in this case we can identify which files and fields need further (manual) processing. For this test task we would like you to create a PoC that
1.) demonstrates the feasibility of the above goal
2.) gets as good accuracy as possible given the time constraints that this is "just" a test

**Materials & Constraints**

You get the following materials from us:
- Training data set of 284 registration certificate scans
- CSV file with filename, license plate, VIN and registration date for each sample in the data set

Since these scans were obtained from actual customers some of the scans are of very similar style and quality, this is a problem you will also find in reality. It is up to you to create a training, validation and test data set from the materials. Please ensure that there is sufficient diversity in your test set, ie. don't use a test set which contains scans from only 2 different customers/styles.

To keep this task simpler we have already separated the registration certificates into individual images for you and have only included german registration certificates in the data set. We hope this helps you focus on the data extraction and accuracy prediction problem.

You are allowed to use any tools of your choosing, including cloud machine learning APIs such as Google OCR or Amazon TextExtract. However, if you do use such tools we ask you to keep the data set on private storages in these services and delete all copies of it upon completion of the project.

**Deliverables**
- All code used to create your PoC, preferably with install & run instructions (no need to get fancy, we would just like to try to run it locally)
- Your training, validation and test data set
- Field level accuracy of your best performing model for the test set
- If your model produces it, information on the confidence of the results achieved on the test set (e.g. probability of correctness for each field)
- Any additional materials which you think would be helpful for us to understand how you approached the problem and how you would proceed from here

For evaluating the project we will of course look at the accuracy of your model in processing the test set as well as the accuracy of the confidence prediction. We prefer a slightly less accurate solution with very accurate confidence over a highly accurate, but on individual samples erratic solution (aka we would rather your models knows when it is unsure than a very bold model who happens to get it right on this test set).

We will also look at the code you created (knowing that this is a test project, so we don't expect things like extensive documentation, refactored & cleaned code etc.) and how you approached the problem. In the end we are interested both in the result that you have achieved as well as your path to get there.