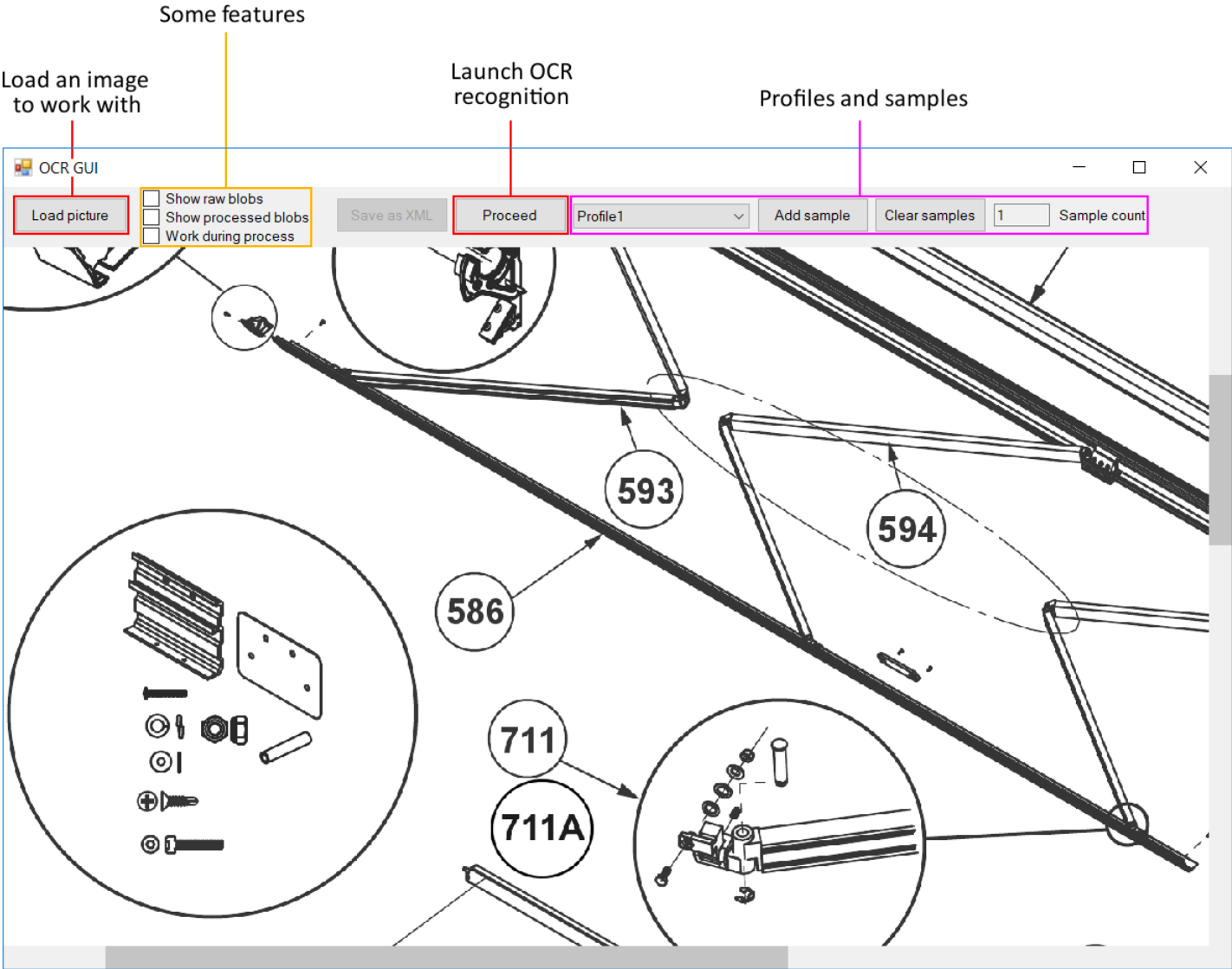


OCR application

Quick user guide

I. Overview

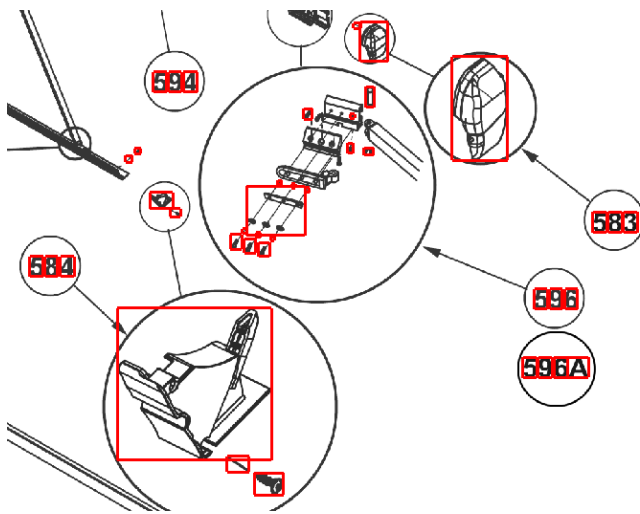


II. Profiles and samples

The goal of the application is to find the regions of interest (the hotspots) before applying the Tesseract OCR recognition. These regions are found using a *profile* that contains *samples*.

Each sample describes the properties (size, ratio, ...) that can be used to find the associated hotspot collection. Some images only need one sample (all the hotspots are the same) whereas some other images require more than one (there are several hotspot collections with different properties).

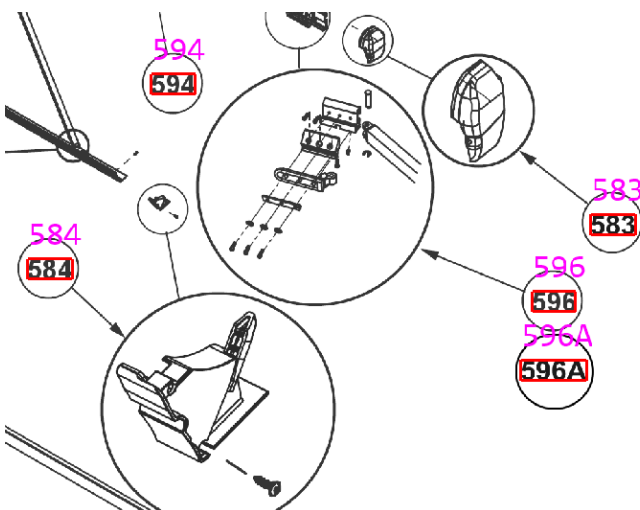
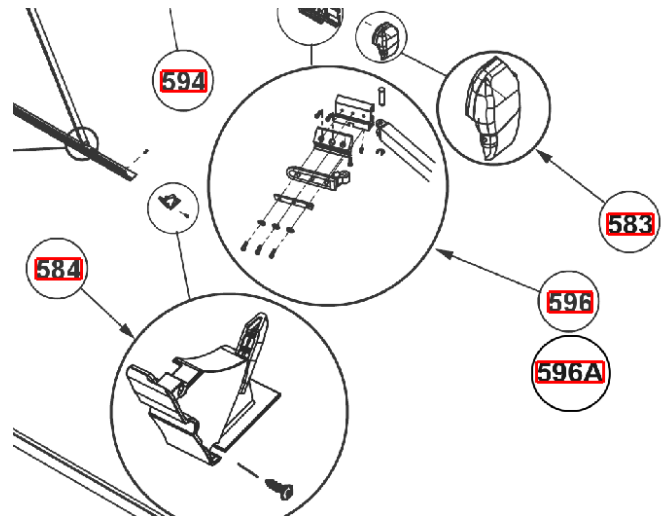
With the checkboxes *Show raw blobs* and *Show processed blobs*, you can visualize all the regions of interest and the ones that have been selected with the samples :



Raw blobs

The very first step of the algorithm

Processed blobs
The hotspots selected with the samples



Final hotspots
The recognized hotspots

Profile config file

You can describe profiles and samples using an XML config file. You need to name it **ocr_profile.config** and place it in **C:\ProgramData\Signifikant\Assert** so that the application can use it (each profile defined can then be selected through the application). The file is reloaded each time an image is loaded in the application.

The file consists of sample and profile declarations :

```
<AllProfiles>
  <Sample name="Sample1">
    ...
  </Sample>

  <Profile name="Profile1">
    <UseSample>Sample1</UseSample>
    ...
  </Profile>
</AllProfiles>
```

You can also declare samples directly in a profile (you don't have to use the tag *UseSample*) :

```
<AllProfiles>
  <Profile name="Profile1">
    <Sample name="Sample1">
      ...
    </Sample>
    ...
  </Profile>

  <Profile name="Profile2">
    <UseSample>Profile1.Sample1</UseSample>
    ...
  </Profile>
</AllProfiles>
```

Profiles and samples can also inherit from other profiles and samples :

```
<Sample name="Sample1">
  ...
</Sample>

<Sample name="Sample2" inherit="Sample1">
  ...
</Sample>

<Profile name="Profile1">
  ...
</Profile>

<Profile name="Profile2" inherit="Profile1">
  ...
</Profile>
```

The XML fields for a profile are :

```
<Profile name="Profile1">
  <UseSample>Sample1</UseSample>

  <Invert>                false      </Invert>
  <EnhanceContrast>       true       </EnhanceContrast>
  <SmoothLevel>           0         </SmoothLevel>
  <BoldLevel>              0         </BoldLevel>
</Profile>
```

- **Invert** is used to invert the image color (used when the image has a black background). Its default value is false.
- **EnhanceContrast** accentuates the contrast of the image. Its default value is true.
- **SmoothLevel** applies a smooth filter on the image. A higher value means a bigger smooth. Its default value is 0, and having a smooth level greater than 1 is uncommon.
- **BoldLevel** applies a bold filter on the image. A higher value means a bigger bold. Its default value is 0, and having a bold level greater than 1 is uncommon.

The smooth filter needs to be used when the image has sharp edges that goes directly from black to white (no gray intermediate pixels between). The bold filter needs to be used when the hotspots are too thin.

The XML fields for a sample are :

```
<Sample name="Sample1">
<!--> <SizeWidth023456789>          19          </SizeWidth023456789>
      <SizeHeight023456789>         30          </SizeHeight023456789>
      <SizeErrorWidth023456789>     0.15        </SizeErrorWidth023456789>
      <SizeErrorHeight023456789>    0.05        </SizeErrorHeight023456789>
      <SizeWidth1>                   13          </SizeWidth1>
      <SizeHeight1>                   30          </SizeHeight1>
      <SizeErrorWidth1>               0.15        </SizeErrorWidth1>
      <SizeErrorHeight1>              0.05        </SizeErrorHeight1>
<!--> <Ratio023456789>              1.55        </Ratio023456789>
      <Ratio1>                        2.16        </Ratio1>
      <RatioError023456789>          0.10        </RatioError023456789>
      <RatioError1>                   0.8         </RatioError1>
<!--> <AlignXSpacingRatio>           1           </AlignXSpacingRatio>
      <AlignYSpacingRatio>            0.10        </AlignYSpacingRatio>
      <AlignHeightRatio>              0.10        </AlignHeightRatio>
<!--> <CircleRadius>                 0           </CircleRadius>
      <CircleError>                   0.5         </CircleError>
<!--> <SmoothLevel>                   0           </SmoothLevel>
      <BoldLevel>                      0           </BoldLevel>
      <HotspotRegex>                   [0-9]+.?    </HotspotRegex>
</Sample>
```

- **SizeWidth023456789** indicates the width (in pixels) of all the digits except 1.
- **SizeHeight023456789** indicates the height of all the digits except 1.
- **SizeWidth1** indicates the width of the digit 1.
- **SizeHeight1** indicates the height of the digit 1.
- **SizeErrorWidth<...>** indicates the accepted error for the width (ratio of the width). Its default value is 0.15 (for both 023456789 and 1).
- **SizeErrorHeight<...>** indicates the accepted error for the height (ratio of the height). Its default value is 0.05 (for both 023456789 and 1).

To make the size property work, you need to specify the width and the height of both *023456789* and *1* digits, or only one of them if you specified some ratio properties (the missing information can then be calculated).

- **Ratio023456789** indicates the ratio height/width of all the digits except *1*. Its default value is 1.5833.
- **Ratio1** indicates the ratio height/width of the digit *1*. Its default value is 2.31.
- **RatioError<...>** indicates the accepted error for the ratio. Its default value is 0.15 (for both *023456789* and *1*).

The ratio property always works since it has default values. If a working size property is defined, the ratio property is no longer used by the application (the size property is more specific).

- **AlignXSpacingRatio** indicates the maximum horizontal space between two blobs to consider them part of the same word/number. If no size property is defined, it is a ratio of the most left blob width. If the size property is defined, it is a ratio of the *SizeWidth023456789* field. Its default value is 1.
- **AlignYSpacingRatio** indicates the maximum vertical space between the center of two blobs to consider them part of the same word/number. It is a ratio of the most left blob height. Its default value is 0.333.
- **AlignHeightRatio** indicates the maximum height of the most right blob to consider it part of the same word/number. It is a ratio of the most left blob height. Its default value is 0.333.

The align property always works since it has default values. Modifying it is very rare, the default values work well in most of the cases.

- **CircleRadius** indicates the radius (in pixel) of the circle around each hotspot.
- **CircleError** indicates the accepted error for the radius (ratio of the radius). Its default value is 0.50.
- **CircleHoughParam2** indicates the second parameter in the Hough Circle Transform algorithm. Its default value is 30.

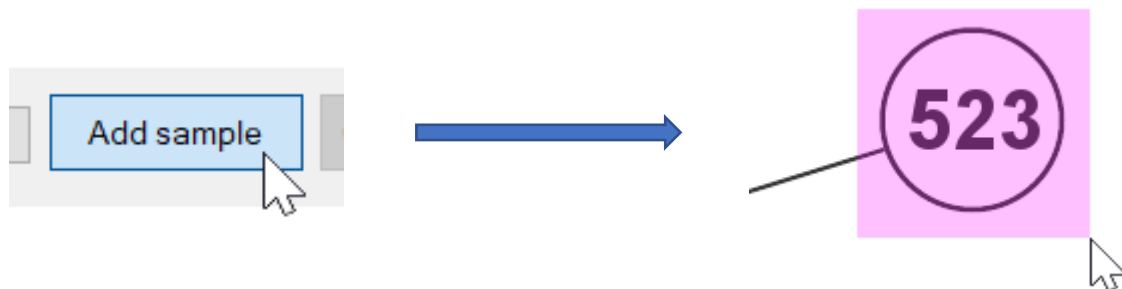
The circle property only works if the radius is greater than 0.5. If the hotspots don't have any circle around them, set the radius to 0 or omit the radius field. The *houghparam2* is a very underground value that indicates the number of votes needed in the algorithm to accept circles. The default value has been empirically found and works in most of the cases.

- **SmoothLevel** applies a smooth filter on the hotspot just before Tesseract recognition. A higher value means a bigger smooth. Its default value is 0, and having a smooth level greater than 1 is uncommon.
- **BoldLevel** applies a bold filter on the hotspot just before Tesseract recognition. A higher value means a bigger bold. Its default value is 0, and having a bold level greater than 1 is uncommon.
- **HotspotRegex** indicates the regular expression the hotspots need to match. Its default value is *[0-9]+.?* (a number that can be followed by one letter).

The smooth and bold level in the samples should be often used instead of the ones in the profiles, since it is applied just before Tesseract recognition, preserving the blob properties.

Sampling from the application

It is possible to retrieve samples directly from the application using the *Add sample* button. Once you click the button, you then need to right click and drag a rectangle around a hotspot :



A new profile named *Sample* is then created, containing one sample. Once the profile is created, using the *Add sample* button again will add new samples in the profile. You can clear the *Sample* profile by using the *Clear samples* button.

The fields used in the profile are the same as the profile that was selected during the sampling process (except for the samples, that are replaced with the one found).

In the retrieved samples, the regular expression and the bold level are used at their default value. Everything else is retrieved directly from the given samples.

Automatic sampling

Automatic sampling is not implemented yet (it should be done for January). Its purpose is to let the application find the samples (based on their statistic properties), and will allow to automatically process several images/folders.

III. Inputs/shortcuts

- **Left click** select/unselect a hotspot
- **Left click (pressed) and drag** select hotspots in the region
- **Shift + left click (pressed) and drag** select/unselect hotspots in the region
- **Right click (pressed) and drag** select the region to process, or select a sample
- **Wheel scroll** scroll vertically
- **Shift + wheel scroll** scroll horizontally
- **Ctrl + wheel scroll** zoom in/out
- **Wheel button (pressed) and drag** free move
- **Double wheel button** center camera on cursor
- **Space** launch the OCR recognition (equivalent to the *Process* button)
- **S** start sampling (equivalent to the *Add sample* button)
- **H (pressed)** hide the hotspots
- **Del** delete the current hotspot selection
- **I** invert the current hotspot selection
- **Ctrl-Z** undo the hotspot selection
- **Ctrl-Y** redo the hotspot selection