

Text Extraction Pipeline Documentation

Author: Falk Böschen

Version 0.6 – July 2016

This document describes the different modules of the pipeline and especially their inputs and outputs. The whole pipeline is based on the streams framework¹. Thus, the communication between different modules is achieved via a Data-Item which is handed from one module to the next. This Data-Item contains all the data and stores them as key-value pairs. The key is a string that is used to index the associated data object which can be of any serializable type. The following module descriptions include the key-values pairs that are used by each module. When a data object that is generated by one module should be used by another one, then the key of the input of the second module must equal the key of the output of the first module. An example can be seen below:

```
<de.cau.cs.kd.processor.open.OpenImage _fileKey="file" imageKey_="image" />
<de.cau.cs.kd.processor.image.grey.ATSLuminanceConversion _imageKey="image" imageKey_="imageGrey" />
<de.cau.cs.kd.processor.image.binary.OtsuHierarchical _imageKey="imageGrey" imageKey_="imageBinary" />
```

The key of the output (`imageKey_`) is equal to the key that is input to the next module (`_imageKey`). If the same key is used for different items, then the old item gets replaced. One can see the parameter naming convention that was used in the example above. The incoming parameters begin with an underscore while the outgoing parameters end with an underscore.

A complete configuration file of the pipeline has to have the following structure:

```
<container>
<stream id="data" class="XXX" ... />

  <process input="data">

    <... />

    <... />

    <... />

  </process>
</container>
```

The whole process is wrapped by a `<container>` tag. First an input has to be specified by defining a stream, which class is one of the source module described below. It is followed by a `<process>` tag which has to have an input key that is equal to the id key of the stream module. Inside the `<process>` tag one can list the processing and sink modules which shall be used.

¹ Streams Framework: <https://sfb876.de/streams/>

| | |
|---------------------------------------|-----------|
| SOURCE MODULES | 4 |
| GUIFILECHOOSER | 4 |
| LOCALFILECHOOSER | 4 |
| PROCESSING MODULES | 5 |
| OPEN FILE MODULES | 5 |
| OPENIMAGE | 5 |
| OPENPDF | 5 |
| IMAGE CONVERSION MODULES | 6 |
| COLORQUANTIZATION | 6 |
| HSVTORGB | 6 |
| RGBTOHSV | 6 |
| ATSC LUMINANCE CONVERSION | 7 |
| NTSC LUMINANCE CONVERSION | 7 |
| EQUAL CONVERSION | 7 |
| CHANNELSELECT | 7 |
| IMAGE BINARIZATION MODULES | 8 |
| THRESHOLDFIX | 8 |
| THRESHOLD DYNAMIC | 8 |
| HISTOGRAM | 8 |
| NIBLACK | 9 |
| OTSU | 9 |
| OTSUCHECKBOARD | 9 |
| OTSUSLIDING | 10 |
| OTSUHIERARCHICAL | 10 |
| OTHER IMAGE PROCESSING MODULES | 10 |
| COLORTOBINARYIMAGES | 10 |
| EDGE | 11 |
| INVERT | 11 |
| MEAN | 11 |
| MEDIAN | 12 |
| CONNECTEDCOMPONENTANALYSIS | 12 |
| PIVOTING | 12 |
| IMAGESFROMCLUSTER | 13 |
| FILTERFORFIGURES | 13 |
| FEATURE MODULES | 14 |
| FEATURESFROMCC | 14 |
| COMPUTEANGLE | 14 |
| FILTER MODULES | 15 |
| AREA | 15 |
| DOTS | 15 |

| | |
|-------------------------------------|-----------|
| SIZE | 16 |
| CLUSTERING MODULES | 16 |
| DBSCAN | 16 |
| GRAVITYGROUPING | 17 |
| GROUPING | 17 |
| MORPHCLUSTERING | 17 |
| MST | 18 |
| OMST | 18 |
| OMSTSUBCLUSTERING | 18 |
| SINGLE | 19 |
| SINGLETON | 19 |
| OCR MODULES | 20 |
| OCROPY | 20 |
| TESSERACT | 20 |
| TEXT PROCESSING MODULES | 21 |
| NORMALIZEWHITESPACES | 21 |
| QUANTITATIVEASSESSMENT | 21 |
| REMOVESPECIALCHARACTERS | 21 |
| REMOVESTRINGS | 21 |
| STRUCTURAL AND OTHER MODULES | 22 |
| CONDENSE | 22 |
| PARALLELIZE | 22 |
| IMAGEVIEW | 22 |
| <u>SINK MODULES</u> | 23 |
| CONSOLE | 23 |
| CREATECSV | 23 |
| CREATEJSON | 23 |

Source Modules

Source modules are used to generate new Data-Items to send along the pipeline. They are defined inside the `class` parameter of a `<stream>` tag. The following modules are all in the `de.cau.cs.kd.source` package and are used to generate Data-Items that contain a file object pointing to a file on the local filesystem.

GUIFileChooser

The **GUIFileChooser** module presents the user with a file chooser dialog which is used to open one file at a time for processing.

| Parameter | Type | Description |
|-----------------------|---------|--|
| <code>id</code> | String | The ID of the stream which is used by <code><process></code> (default: <code>"data"</code>) |
| <code>limit</code> | Long | The maximum number of files to read (default: <code>"-1"</code> – unlimited) |
| <code>fileKey_</code> | String | The key under which the File object is placed (default: <code>"file"</code>) |
| <code>_type</code> | String | The type of files the file chooser can open (default: <code>"."</code> - all files) |
| <code>_log</code> | Boolean | Activate logging functionality (default: <code>"false"</code> - deactivated) |
| <code>_logfile</code> | String | Path to where the log information will be stored (default: null) |

LocalFileChooser

The **LocalFileChooser** module iterates over all files inside a directory and sends one file after the other through the pipeline if the file has the correct type. Recursive mode is possible.

| Parameter | Type | Description |
|-------------------------|---------|--|
| <code>id</code> | String | The ID of the stream which is used by <code><process></code> (default: <code>"data"</code>) |
| <code>limit</code> | Long | The maximum number of files to read (default: <code>"-1"</code> – unlimited) |
| <code>fileKey_</code> | String | The key under which the File object is placed (default: <code>"file"</code>) |
| <code>_type</code> | String | The type of files that will be accepted (default: <code>"."</code> - all files) |
| <code>_input</code> | String | Path to the directory where the files are (default: not specified) |
| <code>_recursive</code> | Boolean | Activate recursive crawling of directory (default: <code>"false"</code> - deactivated) |
| <code>_log</code> | Boolean | Activate logging functionality (default: <code>"false"</code> - deactivated) |
| <code>_logfile</code> | String | Path to where the log information will be stored (default: null) |

Processing Modules

Open File Modules

The following modules are located in the [de.cau.cs.kd.processor.open](#) package.

OpenImage

The **OpenImage** module expects as input a File object which points to an image file that is loaded and stored in the Data-Item.

| Parameter | Type | Description |
|------------------------|---------|--|
| <code>_fileKey</code> | String | The key under which the File object is stored (default: "file") |
| <code>imageKey_</code> | String | The key under which the image is placed as a Double[][][] array (default: "image") |
| <code>_gzip</code> | Boolean | Expect a {ppm pgm pbm}.gz file as input (default: "false") |

OpenPDF

The **OpenPDF** module expects as input a File object which points to a PDF file and extracts all text and images from this PDF and stores it in the Data-Item.

| Parameter | Type | Description |
|-------------------------|--------|--|
| <code>_fileKey</code> | String | The key under which the File object is stored (default: "file") |
| <code>imagesKey_</code> | String | The key under which all images are placed as a List<Double[][][]> (default: not specified) |
| <code>textKey_</code> | String | The key under which the text of the PDF is stored (default: not specified) |

Image Conversion Modules

The following modules are used to convert images into different color codings. A `Double[][][]` structure is used as image format, where the first array dimension is the width of the image, the second dimension is the height of the image, and the third dimension is the color channel.

ColorQuantization

This module reduces the number of colors in an image to a maximal predefined number of colors. The **ColorQuantization** module is located in the `de.cau.cs.kd.processor.image.color` package.

| Parameter | Type | Description |
|--------------------------|---------|--|
| <code>_imageKey</code> | String | The key to access the image (<code>Double[][][]</code>) (default: "image") |
| <code>imageKey_</code> | String | The key under which the modified image is placed as a <code>Double[][][]</code> array (default: "image") |
| <code>_num</code> | Integer | The (maximal) number of colors to reduce to (default: "8") |
| <code>_windowSize</code> | Integer | Parameter for the mean shift algorithm, which defines the pixel neighborhood for computing the mean value (default: "7") |
| <code>_maxDepth</code> | Integer | Parameter for the median cut algorithm, that defines the maximal depth of the 3D kd tree (default: "10") |
| <code>_threshold</code> | Double | Parameter for the mean shift algorithm to define the distance between different colors (default: "25.0") |

HSVtoRGB

This module takes the color channel of an image and converts the values into RGB space, based on the premise that the input color channel had HSV encoding.

The **HSVtoRGB** module is located in the `de.cau.cs.kd.processor.image.color` package.

| Parameter | Type | Description |
|------------------------|--------|--|
| <code>_imageKey</code> | String | The key to access the image (<code>Double[][][]</code>) (default: "image") |
| <code>imageKey_</code> | String | The key under which the modified image is placed as a <code>Double[][][]</code> array (default: "image") |

RGBtoHSV

This module takes the color channel of an image and converts the values into HSV space, based on the premise that the input color channel had RGB encoding.

The **RGBtoHSV** module is located in the `de.cau.cs.kd.processor.image.color` package.

| Parameter | Type | Description |
|------------------------|--------|--|
| <code>_imageKey</code> | String | The key to access the image (<code>Double[][][]</code>) (default: "image") |
| <code>imageKey_</code> | String | The key under which the modified image is placed as a <code>Double[][][]</code> array (default: "image") |

ATSLuminanceConversion

This module converts a color image into a grey image and thus reducing the color channel dimension from 3 to 1. It uses the luminance formula: $y = 0.2126R + 0.7152G + 0.0722B$.

The module is located in the [de.cau.cs.kd.processor.image.grey](#) package.

| Parameter | Type | Description |
|------------------------|--------|---|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>imageKey_</code> | String | The key under which the modified image is placed as a Double[][][] array (default: "image") |

NTSLuminanceConversion

This module converts a color image into a grey image and thus reducing the color channel dimension from 3 to 1. It uses the luminance formula: $y = 0.299R + 0.587G + 0.114B$.

The module is located in the [de.cau.cs.kd.processor.image.grey](#) package.

| Parameter | Type | Description |
|------------------------|--------|---|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>imageKey_</code> | String | The key under which the modified image is placed as a Double[][][] array (default: "image") |

EqualConversion

This module converts a color image into a grey image and thus reducing the color channel dimension from 3 to 1. It uses the luminance formula: $y = (R + G + B) / 3$.

The module is located in the [de.cau.cs.kd.processor.image.grey](#) package.

| Parameter | Type | Description |
|------------------------|--------|---|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>imageKey_</code> | String | The key under which the modified image is placed as a Double[][][] array (default: "image") |

ChannelSelect

This module converts a color image into a grey image and thus reducing the color channel dimension from 3 to 1. It selects the specified color channel and discards the other channels.

The module is located in the [de.cau.cs.kd.processor.image.grey](#) package.

| Parameter | Type | Description |
|------------------------|---------|---|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>imageKey_</code> | String | The key under which the modified image is placed as a Double[][][] array (default: "image") |
| <code>_channel</code> | Integer | The channel to keep. (default: "0") |

Image Binarization Modules

The following modules are located in the `de.cau.cs.kd.processor.image.binary` package.

ThresholdFix

This module converts a grey image into a binary image.

It uses a fixed threshold to decide whether a pixel is assigned with the value 0 or 1.

| Parameter | Type | Description |
|-------------------------|--------|---|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>imageKey_</code> | String | The key under which the modified image is placed as a Double[][][] array (default: "image") |
| <code>_threshold</code> | Double | The threshold for binarization. Range: 0-255 (default: "128.0") |

ThresholdDynamic

This module converts a grey image into a binary image.

It uses a dynamic threshold to decide whether a pixel is assigned with the value 0 or 1.

| Parameter | Type | Description |
|--------------------------|---------|---|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>imageKey_</code> | String | The key under which the modified image is placed as a Double[][][] array (default: "image") |
| <code>_threshold</code> | Double | Maximal difference from the average in the local window (default: "5.0") |
| <code>_windowSize</code> | Integer | Window size (default: "3") |

Histogram

This module converts a grey image into a binary image.

It computes the histogram on the grey values, finds the maximum and sets all values that are above the most common color (minus a threshold) to 1 or else to 0.

| Parameter | Type | Description |
|-------------------------|--------|---|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>imageKey_</code> | String | The key under which the modified image is placed as a Double[][][] array (default: "image") |
| <code>_threshold</code> | Double | Threshold parameter (default: "20.0") |

Niblack

This module converts a grey image into a binary image using Niblack's binarization method.

| Parameter | Type | Description |
|--------------------------|---------|---|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>imageKey_</code> | String | The key under which the modified image is placed as a Double[][][] array (default: "image") |
| <code>_variant</code> | String | The variant of Niblack's method to use. Options: "default", "wolf", "sauvola", "khurshid", "feng" |
| <code>_windowSize</code> | Integer | Parameter defining the neighborhood (default: "42") |
| <code>_k</code> | Double | Parameter used by all methods except "feng" |
| <code>_r</code> | Double | Parameter used by "sauvola", "wolf", and "feng" |
| <code>_a1</code> | Double | Parameter used by "feng" (default: "0.15") |
| <code>_k1</code> | Double | Parameter used by "feng" (default: "0.20") |
| <code>_k2</code> | Double | Parameter used by "feng" (default: "0.03") |

Otsu

This module converts a grey image into a binary image using Otsu's method.

| Parameter | Type | Description |
|------------------------|--------|---|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>imageKey_</code> | String | The key under which the modified image is placed as a Double[][][] array (default: "image") |

OtsuCheckboard

This module converts a grey image into a binary image by subdividing the image in a grid like fashion and applying Otsu's method to each grid cell.

| Parameter | Type | Description |
|-------------------------|---------|---|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>imageKey_</code> | String | The key under which the modified image is placed as a Double[][][] array (default: "image") |
| <code>_gridSizeX</code> | Integer | Number of grid cells in horizontal direction (default: "5") |
| <code>_gridSizeY</code> | Integer | Number of grid cells in vertical direction (default: "5") |

OtsuSliding

This module converts a grey image into a binary image by sliding a window over the image and computing Otsu's method in this window for each pixel.

| Parameter | Type | Description |
|---------------------------|---------|---|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>imageKey_</code> | String | The key under which the modified image is placed as a Double[][][] array (default: "image") |
| <code>_windowSizeX</code> | Integer | Width of the sliding window |
| <code>_windowSizeY</code> | Integer | Height of the sliding window |

OtsuHierarchical

This module converts a grey image into a binary image. It hierarchically subdivides the image where necessary and computes multiple thresholds using Otsu's method to binarize using the edge image.

| Parameter | Type | Description |
|---|---------|---|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>imageKey_</code> | String | The key under which the modified image is placed as a Double[][][] array (default: "image") |
| <code>_edgeBinarizationThreshold</code> | Double | Threshold for binarizing the edge image (default: "32") |
| <code>_minDimensionSize</code> | Integer | Minimum Width/Height of a region to limit the subdivision (default: "15") |
| <code>_hausdorffThreshold</code> | Integer | Threshold for deciding whether to subdivide or not based on the computation of Hausdorff (default: "3") |

Other Image Processing Modules

ColorToBinaryImages

This module takes a color image as input and creates a list of binary images – one for each color that is present in the image, where only the pixels are set to 1 where the original image has that color.

The module is located in the `de.cau.cs.kd.processor.image.binary` package.

| Parameter | Type | Description |
|-------------------------------|--------|--|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>binaryImagesKey_</code> | String | The key to put the List<Double[][][]> of binary images (default: "binaryImages") |

Edge

Computes the edge image of an image.

The module is located in the [de.cau.cs.kd.processor.image.operators](#) package.

| Parameter | Type | Description |
|-----------------------------|--------|---|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>imageKey_</code> | String | The key under which the edge image is placed as a Double[][][] array (default: "image") |
| <code>angleImageKey_</code> | String | The key under which the angle image is placed as a Double[][][] array (default: "imageAngle") |
| <code>gxImageKey_</code> | String | The key under which the horizontal edge image is placed as a Double[][][] array (default: "Gx") |
| <code>gyImageKey_</code> | String | The key under which the vertical edge image is placed as a Double[][][] array (default: "Gy") |
| <code>_method</code> | String | The method to compute the edge image. Options: "sobel" (default), "roberts", "laplace1", and "laplace2" |

Invert

This module inverts the color values of an RGB image.

The module is located in the [de.cau.cs.kd.processor.image.operators](#) package.

| Parameter | Type | Description |
|------------------------|--------|---|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>imageKey_</code> | String | The key under which the modified image is placed as a Double[][][] array (default: "image") |

Mean

This module filters an image with a mean filter, thus replacing each pixel with the average/mean value computed in a local window.

The module is located in the [de.cau.cs.kd.processor.image.operators](#) package.

| Parameter | Type | Description |
|--------------------------|---------|---|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>imageKey_</code> | String | The key under which the modified image is placed as a Double[][][] array (default: "image") |
| <code>_windowSize</code> | Integer | Size of the window (default: "3") |

Median

This module filters an image with a mean filter, thus replacing each pixel with the median value in a local window. The module is located in the [de.cau.cs.kd.processor.image.operators](#) package.

| Parameter | Type | Description |
|--------------------------|---------|---|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>imageKey_</code> | String | The key under which the modified image is placed as a Double[][][] array (default: "image") |
| <code>_windowSize</code> | Integer | Size of the window (default: "3") |

ConnectedComponentAnalysis

This module performs a connected component analysis on a binary image and outputs a list of regions, where each region is a list of pixel coordinates.

The module is located in the [de.cau.cs.kd.processor.image.Label](#) package.

| Parameter | Type | Description |
|--------------------------|--------|---|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>regionsKey_</code> | String | The key for the list<list<pair<int,int>>> of regions (default: "regions") |

Pivoting

This module creates a list of regions by applying a pivoting algorithm on the edge image which alternately computes the histogram projection of the x or y dimension and splitting the image based on the threshold value at the minimal points.

| Parameter | Type | Description |
|----------------------------|---------|---|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>_edgeImageKey</code> | String | The key to access the edge image (Double[][][]) (default: "imageEdge") |
| <code>regionsKey_</code> | String | The key for the list<list<pair<int,int>>> of regions (default: "regions") |
| <code>_threshold</code> | Double | Threshold for subdividing (default : "20.0") |
| <code>_maxDepth</code> | Integer | The maximal subdivision depth (default: unlimited) |

ImagesFromCluster

This module cuts out smaller images from an image which are defined by clusters of regions and stores them as individual Data-Items in a list of Data-Items.

| Parameter | Type | Description |
|--------------------------|---------|--|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>itemsKey_</code> | String | The key under which the list of Data-Items is placed (default: "images") |
| <code>imageKey_</code> | String | The key under which the subimages are placed in their Data-Items (default: "subimage") |
| <code>_clusterKey</code> | String | The key for the cluster information list<list<int>> (default: "cluster") |
| <code>_regionsKey</code> | String | The key to the list<list<pair<int,int>>> of regions (default: "regions") |
| <code>_anglesKey</code> | String | The key for the angles list<Double> (default: "angles") |
| <code>_border</code> | Integer | Parameter defining the width of a border to add to the generated images (default: "5" pixel) |

FilterForFigures

Filters a list of images based on the specified parameters. This module is located in the package `de.cau.cs.kd.processor.image.classify`.

| Parameter | Type | Description |
|----------------------------|---------|--|
| <code>_imageListKey</code> | String | The key under which the list of images (List<Double[][]>) is found (default: "imagelist") |
| <code>imageListKey_</code> | String | The key under which the filtered list of images (List<Double[][]>) is placed (default: "imagelist") |
| <code>_minWidth</code> | Integer | The minimum width in pixel an image must have (default: "500") |
| <code>_minHeight</code> | Integer | The minimum height in pixel an image must have (default: "500") |
| <code>_maxWidth</code> | Integer | The maximum width in pixel an image must have (default: "2000") |
| <code>_maxHeight</code> | Integer | The maximum height in pixel an image must have (default: "2000") |
| <code>_minNumChar</code> | Integer | The min. number of characters that have to be recognized using simple Tesseract OCR (default: "400") |
| <code>_filterWidth</code> | Boolean | Filter images by width if true (default: "true") |
| <code>_filterHeight</code> | Boolean | Filter images by height if true (default: "true") |
| <code>_filterOCR</code> | Boolean | Filter images by <code>_minNumChar</code> if true (default: "true") |

Feature Modules

These modules compute numerical values based on other information.

They are located inside the `de.cau.cs.kd.processor.features` package.

FeaturesFromCC

Computes a list of feature vectors where each entry corresponds to a region and each feature vector consists of width and height of the region, x and y coordinate of its center of mass and its area occupation ratio.

| Parameter | Type | Description |
|---------------------------|--------|--|
| <code>_regionsKey</code> | String | The key to the list<list<pair<int,int>>> of regions (default: "regions") |
| <code>featuresKey_</code> | String | The key to the list<Vector> (default: "features") |

ComputeAngle

Computes the orientation of text lines, which are defined as clusters of regions which are associated with features. Three different methods are available.

| Parameter | Type | Description |
|---------------------------|--------|---|
| <code>_clusterKey</code> | String | The key for the cluster information list<list<int>> (default: "cluster") |
| <code>_featuresKey</code> | String | The key for the features list<Vector> (default: "features") |
| <code>_regionsKey</code> | String | The key for the regions list<list<pair<int,int>>> (default: "regions") |
| <code>anglesKey_</code> | String | The key for the computed angles list<Double> (default: "angles") |
| <code>_method</code> | String | The method to compute the orientation. Options are: "hough", "ssod", and "psd" (default: "hough") |

Filter Modules

The **Filter Modules** are used to filter a list of regions `list<list<pair<int,int>>>` with their features `list<Vector>` using certain thresholds. They are located in the `de.cau.cs.kd.processor.filter` package.

Area

This module filters regions based on the area (of their bounding box) they occupy.

| Parameter | Type | Description |
|---------------------------|---------|---|
| <code>_featuresKey</code> | String | The key for the features <code>list<Vector></code> (default: "features") |
| <code>_regionsKey</code> | String | The key for the regions <code>list<list<pair<int,int>>></code> (default: "regions") |
| <code>featuresKey_</code> | String | The key for the modified features <code>list<Vector></code> (default: "features") |
| <code>regionsKey_</code> | String | The key for the modified regions <code>list<list<pair<int,int>>></code> (default: "regions") |
| <code>_upperratio</code> | Double | Upper limit for the area occupation ratio (default: "1.0" – max) |
| <code>_lowerratio</code> | Double | Lower limit for the area occupation ratio (default: "0.0" – min) |
| <code>_min</code> | Integer | Minimum number of pixel for the region |
| <code>_max</code> | Integer | Maximum number of pixel for the region |

Dots

This module assigns small regions that are most likely dots to the closest region that is not a dot.

| Parameter | Type | Description |
|---------------------------|---------|---|
| <code>_featuresKey</code> | String | The key for the features <code>list<Vector></code> (default: "features") |
| <code>_regionsKey</code> | String | The key for the regions <code>list<list<pair<int,int>>></code> (default: "regions") |
| <code>featuresKey_</code> | String | The key for the modified features <code>list<Vector></code> (default: "features") |
| <code>regionsKey_</code> | String | The key for the modified regions <code>list<list<pair<int,int>>></code> (default: "regions") |
| <code>_merge</code> | Boolean | If set to false, small regions are removed and not merged (default: "true") |
| <code>_size</code> | Integer | Maximum size(width or height) of a dot in pixel (default: "15") |

Size

This module filters regions based on their width and/or height.

| Parameter | Type | Description |
|----------------------------|--------|---|
| <code>_featuresKey</code> | String | The key for the features list<Vector> (default: "features") |
| <code>_regionsKey</code> | String | The key for the regions list<list<pair<int,int>>> (default: "regions") |
| <code>featuresKey_</code> | String | The key for the modified features list<Vector> (default: "features") |
| <code>regionsKey_</code> | String | The key for the modified regions list<list<pair<int,int>>> (default: "regions") |
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>_widthfactor</code> | Double | Factor to multiply the standard deviation of the width distribution to define a range around the average width (default: "3.0") |
| <code>_heightfactor</code> | Double | Factor to multiply the standard deviation of the height distribution to define a range around the average height (default: "3.0") |
| <code>_relativeMax</code> | Double | Maximum width/height defined by a relative factor to the image width/height (default: "1.0") |
| <code>_relativeMin</code> | Double | Minimum width/height defined by a relative factor to the image width/height (default: "0.0") |

Clustering Modules

The following modules are located in the [de.cau.cs.kd.processor.clustering](#) package.

DBSCAN

This module creates a clustering on regions based on their feature vectors.

| Parameter | Type | Description |
|-----------------------------|---------|--|
| <code>_featuresKey</code> | String | The key for the features list<Vector> (default: "features") |
| <code>clusterKey_</code> | String | The key for the cluster information list<list<int>> (default: "cluster") |
| <code>_adaptive</code> | String | List of indices of feature vector entries to use for computing the <code>_eps</code> parameter adaptively [optional] |
| <code>_eps</code> | Double | Parameter defining the neighborhood of a region (default: "15.0") |
| <code>_minPts</code> | Integer | The minimum number of points required to start a cluster (default: "1") |
| <code>_separateNoise</code> | Boolean | If "true", a separate cluster for all regions that are classified as noise is created. Otherwise the noise objects are discarded. (default: "false") |

GravityGrouping

This module groups/clusters regions using a method proposed by Weihua Huang et al.² which is based on Newton's Formula for gravity.

| Parameter | Type | Description |
|---------------------------|--------|---|
| <code>_featuresKey</code> | String | The key for the features list<Vector> (default: "features") |
| <code>_regionsKey</code> | String | The key for the regions list<list<pair<int,int>>> (default: "regions") |
| <code>clusterKey_</code> | String | The key for the cluster information list<list<int>> (default: "cluster") |
| <code>_threshold</code> | Double | Threshold for the distance that is computed by the gravity formula. (default: "20.0") |

Grouping

This module groups regions using another method proposed by Weihua Huang and C.L. Tan.³

| Parameter | Type | Description |
|---------------------------|--------|--|
| <code>_featuresKey</code> | String | The key for the features list<Vector> (default: "features") |
| <code>_regionsKey</code> | String | The key for the regions list<list<pair<int,int>>> (default: "regions") |
| <code>clusterKey_</code> | String | The key for the cluster information list<list<int>> (default: "cluster") |
| <code>_k</code> | Double | Parameter in the grouping function (default: "10.0") |

MorphClustering

A morphological clustering based on the work of Chiang and Knoblock.⁴

| Parameter | Type | Description |
|---------------------------------|--------|--|
| <code>_featuresKey</code> | String | The key for the features list<Vector> (default: "features") |
| <code>_regionsKey</code> | String | The key for the regions list<list<pair<int,int>>> (default: "regions") |
| <code>clusterKey_</code> | String | The key for the cluster information list<list<int>> (default: "cluster") |
| <code>_maxCurvatureRatio</code> | Double | Parameter for the morphological clustering (default: "1.2") |
| <code>_maxSizeRatio</code> | Double | Parameter for the morphological clustering (default: "2.0") |
| <code>_maxDistanceRatio</code> | Double | Parameter for the morphological clustering (default: "2.0") |

² Huang, W.; Tan, C. L. & Leow, W. K. Associating Text and Graphics for Scientific Chart Understanding Eighth International Conference on Document Analysis and Recognition (ICDAR 2005), 29 August - 1 September 2005, Seoul, Korea, IEEE Computer Society, 2005, 580-584

³ Huang, W. & Tan, C. L. King, P. R. & Simske, S. J. (Eds.) A system for understanding imaged infographics and its applications Proceedings of the 2007 ACM Symposium on Document Engineering, Winnipeg, Manitoba, Canada, August 28-31, 2007, ACM, 2007, 9-18

⁴ Chiang, Y. & Knoblock, C. A. Recognizing text in raster maps Geoinformatica, 2015, 19, 1-27

MST

This module applies a minimum spanning tree clustering onto the regions. The MST can be split using one of two methods: Threshold-based, or Inconsistency-based.

| Parameter | Type | Description |
|-----------------------------|---------|--|
| <code>_featuresKey</code> | String | The key for the features list<Vector> (default: "features") |
| <code>clusterKey_</code> | String | The key for the cluster information list<list<int>> (default: "cluster") |
| <code>_threshold</code> | Double | A length threshold for removing edges |
| <code>_inconsistency</code> | Boolean | If set to "true", the average edge length at a node is computed and all edges longer than the average times 1.2 are removed. |

OMST

This module subdivides an existing clustering using a minimum spanning tree approach that is split using orientation information.

| Parameter | Type | Description |
|---------------------------|--------|--|
| <code>_featuresKey</code> | String | The key for the features list<Vector> (default: "features") |
| <code>_clusterKey</code> | String | The key for the cluster information list<list<int>> (default: "cluster") |
| <code>clusterKey_</code> | String | The key for the new cluster information list<list<int>> (default: "cluster") |
| <code>_angle</code> | String | Two comma separated angles in degree that define the range around the mean orientation of a cluster where all edges with orientations outside this range are removed. (default: "60,60") |

OMSTSubclustering

A slightly different version of the implementation of the [OMST](#) module.

| Parameter | Type | Description |
|---------------------------|--------|--|
| <code>_featuresKey</code> | String | The key for the features list<Vector> (default: "features") |
| <code>_clusterKey</code> | String | The key for the cluster information list<list<int>> (default: "cluster") |
| <code>clusterKey_</code> | String | The key for the new cluster information list<list<int>> (default: "cluster") |
| <code>_angle</code> | String | Two comma separated angles in degree that define the range around the mean orientation of a cluster where all edges with orientations outside this range are removed. (default: "60,60") |

Single

This module creates a single cluster containing all regions.

| Parameter | Type | Description |
|---------------------------|--------|---|
| <code>_featuresKey</code> | String | The key for the features list<Vector> (default: "features") |
| <code>clusterKey_</code> | String | The key for the cluster information list<list<int>> (default: "cluster") |

Singleton

This module creates an individual cluster for every region.

| Parameter | Type | Description |
|---------------------------|--------|---|
| <code>_featuresKey</code> | String | The key for the features list<Vector> (default: "features") |
| <code>clusterKey_</code> | String | The key for the cluster information list<list<int>> (default: "cluster") |

OCR Modules

The following modules are located in the `de.cau.cs.kd.processor.ocr` package.

Ocropy

This module is accessing the Ocropy installation on a Unix system via command line to recognize the text inside an image and returning the recognized text as a String.

| Parameter | Type | Description |
|-------------------------|--------|---|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>textKey_</code> | String | The key to store the recognized text (default: "ocredtext") |
| <code>_tmpfolder</code> | String | The folder which the Ocropy installation shall use for its temporary data (default: "book") |
| <code>_option</code> | String | The parameters that shall be passed to Ocropy (default: none) |

Tesseract

This module uses the Tess4J API to apply Tesseract OCR on the input image. Output is a text String.

| Parameter | Type | Description |
|------------------------|--------|---|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>textKey_</code> | String | The key to store the recognized text (default: "ocredtext") |
| <code>_language</code> | String | The language model Tesseract shall use (default: "eng") (based on available dataset) |
| <code>_OEMmode</code> | String | Set the OEM mode ⁵ of tesseract (default: "OEM_TESSERACT_ONLY") |
| <code>_PSMmode</code> | String | Set the PSM mode ⁶ of tesseract (default: "PSM_OSD_ONLY") |

⁵ Options are:

"OEM_CUBE_ONLY", "OEM_DEFAULT", "OEM_TESSERACT_ONLY", "OEM_TESSERACT_CUBE_COMBINED".

See <http://tess4j.sourceforge.net/docs/docs-3.0/net/sourceforge/tess4j/ITessAPI.TessOcrEngineMode.html>

⁶ See <http://tess4j.sourceforge.net/docs/docs-3.0/net/sourceforge/tess4j/ITessAPI.TessPageSegMode.html>

Text Processing Modules

These modules can be used to process text strings.

The following modules are located in the `de.cau.cs.kd.processor.text` package.

NormalizeWhitespaces

This module replaces all whitespaces (even sequences) with single blank spaces.

| Parameter | Type | Description |
|-----------------------|--------|---|
| <code>_textKey</code> | String | The key to the String to process (default: <code>"ocredtext"</code>) |
| <code>textKey_</code> | String | The key where the processed String is stored (default: <code>"processedtext"</code>) |

QuantitativeAssessment

Sets a String to null if its number of characters divided by the number of regions that were used to create these characters is greater than the `_maxRatioDifference`.

| Parameter | Type | Description |
|----------------------------------|--------|---|
| <code>_textKey</code> | String | The key to the String to process (default: <code>"ocredtext"</code>) |
| <code>textKey_</code> | String | The key where the processed String is stored (default: <code>"processedtext"</code>) |
| <code>_maxRatioDifference</code> | Double | The maximal ration between characters and regions (default: <code>"0.1"</code>) |

RemoveSpecialCharacters

Removes special characters from a String.

| Parameter | Type | Description |
|-----------------------|--------|---|
| <code>_textKey</code> | String | The key to the String to process (default: <code>"ocredtext"</code>) |
| <code>textKey_</code> | String | The key where the processed String is stored (default: <code>"processedtext"</code>) |
| <code>_remove</code> | String | A string of the characters to removed (default: all special characters) |

RemoveStrings

Sets text strings that contain more than `_num` special characters to null.

| Parameter | Type | Description |
|-----------------------|---------|--|
| <code>_textKey</code> | String | The key to the String to process (default: <code>"ocredtext"</code>) |
| <code>textKey_</code> | String | The key where the processed String is stored (default: <code>"processedtext"</code>) |
| <code>_remove</code> | String | A string of the characters which are detected to decide whether to remove a String (default: all special characters) |
| <code>_num</code> | Integer | The max. allowed number of special characters |

Structural and other Modules

Condense

This module merges a list of multiple Data-Items into one Data-Item.

This module is located inside the [de.cau.cs.kd.processor](#) package.

| Parameter | Type | Description |
|------------------------|---------|---|
| <code>_itemsKey</code> | String | The key to the Data-Item list (default: none) |
| <code>_mergeKey</code> | String | A list of comma separated keys which shall be merged (default: none) |
| <code>_append</code> | Boolean | Merges attributes that are list by extending the list and not creating a list over lists. (default: "true") |
| <code>_keepNull</code> | Boolean | If set to false, null attributes are skipped (default: "true") |

Parallelize

This module splits a Data-Item into multiple Data-Items for parallel processing while splitting an attribute and duplicating the rest. This module is located inside the [de.cau.cs.kd.processor](#) package.

| Parameter | Type | Description |
|------------------------|--------|---|
| <code>_splitKey</code> | String | The key of the list or array attribute to split (default: none) |
| <code>itemsKey_</code> | String | The key where the list of newly created Data-Items are stored (default: none) |

ImageView

This module shows an image with the option to additionally visualize regions or cluster.

This module is located inside the [de.cau.cs.kd.processor.view](#) package.

| Parameter | Type | Description |
|---------------------------|--------|--|
| <code>_imageKey</code> | String | The key to access the image (Double[][][]) (default: "image") |
| <code>_regionsKey</code> | String | The key to access the regions (default: null) [optional] |
| <code>_featuresKey</code> | String | The key to access features (default: null) [optional] (required for regions and cluster visualization) |
| <code>_clusterKey</code> | String | The key to access the clusters (default: null) [optional] |
| <code>_file</code> | String | The path to where the image shall be saved (default: none) [optional] |

Sink Modules

The Sink Modules are designed to output the content of the current Data-Item. They are not real sinks that terminate the execution, but a different kind of processing modules. Thus, they need to be inside the `<process>` tags as well. They are located inside the `de.cau.cs.kd.sink` package.

Console

The **Console** module prints all contents of the Data-Item to the standard output, except for nested lists, which are skipped. It outputs each key together with its data. No parameters are used.

| Parameter | Type | Description |
|-----------|------|-------------|
| - | - | - |

CreateCSV

The **CreateCSV** module generates a file where each line represents an object which is described by the values in the different columns which are separated by a pre-defined symbol. It requires one or multiple lists as input where the i^{th} item in one list describes the same object as the i^{th} item in any of the other lists.

| Parameter | Type | Description |
|-------------------------|--------|---|
| <code>_file</code> | String | Path to where the CSV file will be stored (default: not specified) [optional] (required if <code>_fileKey</code> not set) |
| <code>_fileKey</code> | String | The key under which the File object is located (default: null) [optional] (required if <code>_file</code> not set) |
| <code>_separator</code> | String | The separating character that is used to separate the different columns in the generated file (default: <code>","</code>) [optional] |
| <code>_itemsKey</code> | String | A comma separated list of keys that all point to list object containing the values to output (default: null) |
| <code>_suffix</code> | String | A suffix that is added to the generated files name. (default: none) [optional] |
| <code>_prefix</code> | String | A prefix that is added to the generated files name. (default: none) [optional] |

CreateJSON

The **CreateJSON** module works similar like the **CreateCSV** module but instead of a CSV file, it creates a JSON file with additional meta-information.

| Parameter | Type | Description |
|----------------------------|--------|--|
| <code>_fileKey</code> | String | The key under which the File object is located (default: null) for meta-data extraction |
| <code>_imageKey</code> | String | The key under which the image data is located (default: null) |
| <code>_itemsKey</code> | String | A comma separated list of keys that all point to list object containing the values to output (default: null) |
| <code>_outputFolder</code> | String | Path to the directory where the JSON file(s) will be stored (default: not specified) [optional] |