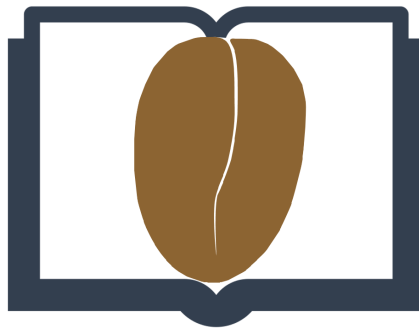

PrepOCReSSor – The QCRI Preprocessing Tool for OCR

Version 0.2

*Qatar Computing Research Institute, HBKU
Felix Stahlberg and Stephan Vogel*

2015-07-24



معهد قطر لبحوث الحوسبة
Qatar Computing Research Institute

عضو في مؤسسة قطر
Member of Qatar Foundation

Abstract

This document describes the capabilities of the open source software *PrepOCReSSor*. The tool is developed at the Qatar Computing Research Institute for preprocessing document images for optical character recognition. The tool follows the pipeline paradigm in Unix-like operating systems: A set of image processing operations is chained such that the output of each operation serves as input to the next one. The tool supports batch processing for high parallelism and scalability. The OpenCV (Bradski and Kaehler, 2008) library provides efficiently implemented computer vision algorithms and a efficient infrastructure. *PrepOCReSSor* is intended to be used in combination with the recognition toolkit *Kaldi* (Povey et al., 2011) and supports file formats for feature sets (.ark,t) and forced-alignments (.al) for a seamless integration. Even though we focus on Arabic script, the tool has been successfully used for other writing systems, e.g. Latin in the ICDAR2015 Competition HTRtS on historic documents.

Contents

1	Introduction	4
1.1	Preprocessing for Optical Character Recognition	4
1.2	PrepOCReSSor Design Philosophy	4
1.3	Licence	5
2	Using PrepOCReSSor	5
2.1	Installation	5
2.2	Getting Help	7
2.3	Tutorials	7
2.4	Example Pipelines	11
3	Troubleshooting	15
4	Operation and Parameter Reference	15
4.1	Global Parameters	15
4.2	Operations	18

1 Introduction

1.1 Preprocessing for Optical Character Recognition

Offline optical character recognition (OCR) refers to the conversion of printed, typewritten, or handwritten text in a scanned image to machine-encoded text. State-of-the-art OCR systems are based on Hidden Markov Models (HMMs) which are statistical models for sequences of feature vectors. The order of the feature vectors within the sequence can represent temporal dependencies. For instance, in automatic speech recognition (ASR), the audio recording is often split into 10-15 ms chunks and a feature vector is extracted for each of these chunks (Huang et al., 2001) (Fig. 1(a)). The input for OCR, however, are images of single text lines. Analogously to ASR, we split the image into chunks with three pixel width and extract a feature vector for each of these chunks. The sequential order of the extracted feature vectors is defined by the reading order of the script (e.g. right-to-left for Arabic in Fig. 1(b)).

Realizing the similarities between offline OCR and ASR, we suggest to apply the state-of-the-art speech recognition toolkit *Kaldi* (Povey et al., 2011) for recognizing text in scanned documents. The *PrepOCReSSor* tool bridges the gap between ASR and OCR and prepares images in a way that they can be passed through to Kaldi. *PrepOCReSSor* provides comprehensive functionality to break down the initial document image into text lines, and convert each text line to a sequence of feature vectors for training or decoding with Kaldi. Therefore, the main functions of *PrepOCReSSor* can be grouped into one of the following categories:

- **Document layout analysis:** Document rotation, text/non-text segmentation, line-segmentation etc.
- **Text image normalization:** Baseline estimation, slant correction, pen size normalization, letter size normalization etc.
- **Feature extraction:** Methods for feature vector extraction from normalized text line images.

The overall goal of *PrepOCReSSor* is to provide a comprehensive addition to Kaldi for OCR research. The combination of Kaldi and *PrepOCReSSor* results in a fully-fledged scalable OCR framework with state-of-the-art recognition performance.

1.2 PrepOCReSSor Design Philosophy

We defined the following non-functional requirements as the main design goals for *PrepOCReSSor*.

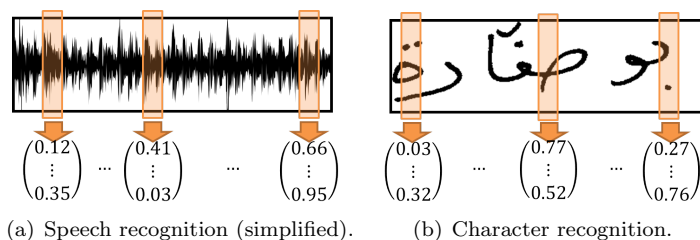


Figure 1: Feature extraction in speech recognition and optical character recognition.

- **Modularity:** In agreement with the Unix philosophy we implement the concept of modularity. PrepOCReSSor provides a large number of small, specialized operations which can be composed in a single Unix-like pipeline. The advantage of this approach is that PrepOCReSSor is highly customizable and can be applied to a wide range of different tasks – i.e. different document types or writing systems. The disadvantage is the manual effort of composing the pipeline for the task at hand. Section 2.4 lists a number of examples which can serve as starting point for your own experiments.
- **Scalability:** PrepOCReSSor supports batch processing. The documents can be distributed to any number of threads. As each document is processed separately, no inter-thread communication is necessary and we achieve nearly a linear (ideal) speed up.
- **Efficiency:** The OpenCV library (Bradski and Kaehler, 2008) provides performance-optimized code for basic computer vision and is widely used both in academia and industry. PrepOCReSSor makes heavy use of algorithms and operations provided by this library.

1.3 Licence

PrepOCReSSor (Copyright ©2015, QCRI a member of Qatar Foundation. All Rights Reserved) is licensed under the Apache License, Version 2.0 (the "License"); you may not use it except in compliance with the License. You may obtain a copy of the License at

<http://www.apache.org/licenses/LICENSE-2.0>.

Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License.

A copy of the License can be found in the LICENSE file in the root directory.

2 Using PrepOCReSSor

2.1 Installation

PrepOCReSSor is implemented in Java and platform-independent. Therefore, it is possible to run it on iOS, Windows, Linux/Mac, and Android.

2.1.1 Installation on Ubuntu Linux

The following instructions explain the installation on Debian-based systems like Ubuntu but can be easily extended for other platforms. The commands in this guide should work in standard Unix shells like zsh and bash. It was tested on Ubuntu 15.04.

1. **Install the Java runtime environment.** PrepOCReSSor was tested with Java 1.7 but should run with other versions as well. On Ubuntu, Java is installed by default. You can check the version number by typing `java -version` into your shell.
2. **Install the OpenCV library.** PrepOCReSSor was tested with OpenCV 2.4.10 but other 2.4.x versions are likely to work. Ubuntu

provides out-of-the-box packages which can be installed with the following command:

```
sudo apt-get install libopencv2.4-java
```

If you are not using Ubuntu, you can download the latest OpenCV 2.4.x version from <http://opencv.org/downloads.html> and follow the installation instructions¹.

3. **Download PrepOCReSSor.** The easiest way to get started with PrepOCReSSor is to download the latest release from `bla` and unzip it wherever you like to install PrepOCReSSor. Alternatively, you can compile PrepOCReSSor by yourself. The repository at <https://bitbucket.org/fstahlberg/preprocessor> contains an Eclipse project including the source code itself as well as the PRiMA library² and Apache Commons Math³.
4. **Configure PrepOCReSSor.** If you don't use Ubuntu Linux or you compiled OpenCV by yourself without using the Ubuntu packages, you need to tell PrepOCReSSor where to find the OpenCV library. Open the `preprocessor` file in the installation root directory in your favourite text editor. You need to set the variables `OPENCV_JAR_PATH` and `OPENCV_NATIVE_LIB`. The variable `OPENCV_JAR_PATH` should point to the OpenCV `.jar` file. For example, in OpenCV 2.4.10 this file can be found within the OpenCV installation in `<opencv-install-dir>/bin/opencv-2410.jar`. If you don't find it, you may have compiled OpenCV without Java support. The `OPENCV_NATIVE_LIB` variable needs to contain the native library directory path (usually `<opencv-install-dir>/lib`). This directory should contain a file called `libopencv_java2410.so` or similar.
5. **Test PrepOCReSSor installation.** You can start PrepOCReSSor by changing into the installation directory and type the following command into the shell:

```
./preprocessor -help
```

This should output a list of global parameters together with a description for each of them. To test if the OpenCV library is installed and configured correctly, type

```
./preprocessor
```

(i.e. without arguments). The output should be similar to this:

```
13:33:14 INFO: Configuration loaded ...
13:33:14 FATAL: Input file 'imageList.txt' reading
error: imageList.txt (No such file or directory)
```

→ Section 3

If you get a significantly different output, consult Section 3 for troubleshooting.

¹For more information about the Java support of OpenCV, check http://docs.opencv.org/doc/tutorials/introduction/desktop_java/java_dev_intro.html

²Original available at <http://primaresearch.org/tools/PAGELibraries>

³Available at <https://commons.apache.org/proper/commons-math/>. PrepOCReSSor was tested with Apache Commons Math 3.4.1

6. **Make your shell aware of PrepOCReSSor.** This manual assumes that you have included PrepOCReSSor in your `$PATH` environment variable so that you can start it with typing `prepocressor` into your shell. You can do this by writing the following line at the end of your `~/ .bashrc`:

```
export PATH=$PATH:<prepocressor-install-dir>
```

Alternatively, you can create a symlink to PrepOCReSSor in a directory which is already in your `$PATH` variable.

```
sudo ln -s /usr/local/bin/prepocressor
      <prepocressor-install-dir>/prepocressor
```

2.1.2 Installation on Other Platforms

As PrepOCReSSor is written in Java, the application is platform-independent and can run on a wide range of operating systems. Modifying the instructions in the previous section for other Linux distributions should be straight-forward. For Windows, however, the `prepocressor` script in the PrepOCReSSor root directory needs to be adjusted to Windows syntax.

2.2 Getting Help

The `-help` parameter in PrepOCReSSor displays all available global parameters together with a short description.

```
prepocressor -help
```

Detailed help texts for specific operations can be displayed by adding the names of the operations. For example, the following command shows detailed descriptions for the `log` and `tee` operations.

```
prepocressor -help log tee
```

For help with modifying and extending the code (Section 2.3.4) you can find the JavaDoc in the `javadoc/` subdirectory.

2.3 Tutorials

2.3.1 Specifying Input and Output Files

The most important parameters for specifying the in- and output of the PrepOCReSSor pipeline are `-inputFile` and `-outputPath`. Following formats are supported:

- Windows bitmaps - `*.bmp`, `*.dib`
- JPEG files - `*.jpeg`, `*.jpg`, `*.jpe`
- JPEG 2000 files - `*.jp2`
- Portable Network Graphics - `*.png`
- Portable image format - `*.pbm`, `*.pgm`, `*.ppm`
- Sun rasters - `*.sr`, `*.ras`
- TIFF files - `*.tiff`, `*.tif`
- Comma-separated values - `*.csv`

The following command loads an image in JPEG format (`test.jpg`) and stores it without modification in PNG format.

```
prepocressor -inputFile test.jpg -outputPath test.png
```

PrepOCReSSor also supports a batch mode which is particularly useful in combination with the `-nThreads` parameter. In the following example, `test.txt` is a text file containing a newline-separated list of paths to image files. PrepOCReSSor converts each of these files to a grayscale image using 8 threads. Per default, the output images are stored with a time and user name encoding prefix.

```
prepocressor -inputFile test.txt
              -nThreads 8 -pipeline "grayscale"
```

In batch mode, a simple `-outputPath` like `"test.png"` is not useful because the results for all images in the batch are written to `test.png` – the same file gets overridden multiple times, and `test.png` just contains the processed image for the last entry in the batch. Therefore, you can use special placeholders in the `-outputPath` parameter:

- `%dir`: Name of the directory containing the input image.
- `%base`: Base name of the input image file name (without file extension).
- `%-base`: Same as `%base`, but cut after first minus hyphn ('-')
- `%idx`: Some operators split up input images into smaller pieces. The pieces are stored subsituting
- `%ext`: File extension (given by the input file).

The following command is similar to the previous example but stores the generated files in a dedicated folder called `blacknwhite`.

```
prepocressor -inputFile test.txt -pipeline "grayscale"
              -outputPath "blacknwhite/%base%ext"
```

2.3.2 Basic Image Manipulation

PrepOCReSSor offers a number of basic image transformation operations which are not necessarily related to OCR or document image processing. They are a great way to get used to the pipeline architecture of PrepOCReSSor.

The following command first stretches the image along the x-axis and then transposes it. This results in an image which is stretched horizontally and then flipped such that the height of the resulting image is two times the width of the original image.

```
prepocressor -inputFile test.png -outputPath out.png
              -pipeline "scale _-xScale_2|transpose"
```

If we switch the order of the `scale` and `transpose` operation, we produce an image which is stretched vertically and turned on its side.

```
prepocressor -inputFile test.png -outputPath out.png
              -pipeline "transpose|scale _-xScale_2"
```

To reproduce the same result as in the first command, we need to stretch around the y-axis after transposing the image:


```

prepocressor -inputFile test.png -outputPath out.png
              -pipeline "transpose | scale -yScale 2"

```

2.3.3 Feature Extraction

→ Section 4.2.19

The `featExtract` operation (Section 4.2.19) is the main interface to the Kaldi toolkit. The input is expected to be a single channel image of a single normalized text line. The `-extractors` parameter specifies which feature extractors are used. If you specify multiple extractors, the features are stacked on each other. The `featExtract` operation accepts a number of extractor specific parameters. By convention, `-featXyz*` parameters are specific to the `xyz` extractor. Following extraction methods are implemented:

- *raw* – Raw pixel values.
- *directional* – Directional features.
- *snake* – Snake feature extraction method (also called segment-based method in (Stahberg and Vogel, 2015))
- *runlengths* – Use pixel-wise runlengths in 4 directions
- *anhdf* – ANHDF features (El-Mahallawy, 2008)
- *distribution* – Distribution features as defined by (Likforman-Sulem et. al., 2012)
- *concavity* – Concavity features as defined by (Likforman-Sulem et. al., 2012)

The following hints help you working with the `featExtract` operation.

→ Section 4.2.22

- PrepOCressor was initially designed for Arabic script with a right-to-left reading direction. Therefore, the `featExtract` method reads the image from right to left to extract the feature vector sequences. If you deal with a left-to-right writing system (e.g. Latin), insert a `flip` operation (Section 4.2.22) prior to `featExtract`.
- The produced feature files are in `.ark,t` format (i.e. text) and therefore very large. You should compress them with Kaldi's `copy-feats` command:

```

copy-feats ark,t:data/feats/pixel_test.ark,t
           ark,scp:data/feats/pixel_test.ark,
           data/test/feats.scp

```

- The feature vector dimension of many feature extractors is dependent on the image height. Therefore, it is important that all images in the pipeline have the same height. Also, keep in mind the curse of dimensionality and that a high dimensionality leads to huge feature files. In our experiments, we use image heights of 40-70 pixels.
- Set `-nThreads` to 1 when using the `featExtract` operation. This ensures that the initial order of images is preserved. Kaldi can be bitchy when it comes to the order of the entries.

→ Section 2.4

You can find some example pipelines for feature extraction in Section 2.4.

2.3.4 Modifying the Code

Instructions for accessing the PrepOCReSSor repository can be found at following URL:

<https://bitbucket.org/fstahlberg/preprocessor>

This section outlines the high level design of the project. It consists of six packages:

- `qa.qcri.preprocessor.datastructures` – This package contains basic data structures like images or image lists that are used to transfer data through the pipeline.
- `qa.qcri.preprocessor.imageprocessing` – This package contains tools for image processing.
- `qa.qcri.preprocessor.io` – This package contains classes for I/O handling, i.e. loading data sets, logging, storing results.
- `qa.qcri.preprocessor.operations` – This package contains all available operations, i.e. all commands which can be used within the pipeline.
- `qa.qcri.preprocessor.operations.feats` – The `feats` package contains feature extractors that work together with the `featExtract` operation
- `qa.qcri.preprocessor.ui` – This package contains classes for the user interaction.

The main runner class is `qa.qcri.preprocessor.ui.Main`. If you want to implement a new operation, you need to inherit from the `Operation` class in the `operations` package. Take a look at the `FlipOperation` class as a basic example. Operations are required to implement at least two methods:

- `createConfiguration()` – This allows you to define the possible parameters for the operation and insert a description. The `flip` operations allows one integer parameter called `flipCode` which decides the axis along which to flip the image. Integer, String, and Float parameters are supported. The type is derived from the type of the second argument of `addParameter` (the default value).
- `processIndividual()` – This method contains the actual implementation of the operation.

Images in the pipeline are represented as `datastructures.Individual` instances. The `Individual` class stores the image itself (see `Individual.getContent()`) together with some meta information. The `processIndividual()` returns a list of individuals because operations can split up images in the pipeline into smaller parts. As `flip` does no such thing, the method returns a list with a single entry holding the input individual. This is possible because the OpenCVs `flip` implementation is implemented in a in-place manner.

Open the `ScaleOperation` class for a more complex example. In this operation, a new image is created (`dst`) and the result is written to that image. The returned list contains a freshly created `Individual` instance referring to `dst`. Note that Java's garbage collector does not apply to OpenCV matrices. Therefore, you need to release `Mat` instances after

usage to prevent memory leaks. Of course, do not release matrices which are passed through the pipeline.

You should store new operation classes in the `operations` package. The naming convention is `XYZOperation`. You can call your operation in the `PrepOCReSSor` pipeline with `xyz`. Please add your new operation to the list in `GlobalConfiguration` in the `ui` package to make it visible in the documentation.

2.4 Example Pipelines

2.4.1 QCRI Submission for the ICDAR2015 Competition HTRtS

This section describes the `PrepOCReSSor` for the QCRI submission to the *ICDAR2015 Competition HTRtS: Handwritten Text Recognition on the tranScriptorium Dataset* Sanchez et al. (2015). For all but the `2ndBatch` set, the following pipeline has been used to generate binarized text line images:

```
invert |
multiChannelOtsu
    -xmlPath <page-dir>/PAGE/%base%idx.xml
    -blackDiscount 0.1
    -normalizeRegionChannels |
morph
    -operation close
    -kernelSize 2
    -kernelShape ellipse |
cutWithPageXml
    -xmlPath <page-dir>/%base%idx.xml
    -extractRegions 0
    -extractTextObjects -usePageIds |
normalize
```

For the `2ndBatch`, no line segmentation was given, so we applied our line segmentation algorithm based on fitting a sinus function to the vertical projection profile.

```
invert |
multiChannelOtsu
    -xmlPath <page-dir>/%base%idx.xml
    -blackDiscount 0.2
    -normalizeRegionChannels
    -maxForegroundFraction 0.1
    -extractRegions 1
    -extractTextObjects 0 |
morph
    -operation close
    -kernelSize 2
    -kernelShape ellipse |
cutWithPageXml
    -xmlPath <page-dir>/%base%idx.xml
    -extractRegions 1
    -extractTextObjects 0 |
projectionLineSegmentation
    -minLineHeight 100
    -maxLineHeight 320
    -analysisMode 0 |
```

```

vertTextSegmentation
  -minWidth 0.0005
  -minMargin 0.05
  -minSlope 0.000001
  -morph closeFirst
  -concatChildren |
normalize

```

The binarized images were then normalized using the following pipeline:

```

grayscale |
transpose |
removeUnderline
  -minRelWidth 0.2
  -minWidth 10
  -foregroundThreshold 20
  -maxVariation 5
  -maxHeight 150 |
flip -flipCode 0 |
removeUnderline
  -minRelWidth 0.2
  -minWidth 10
  -foregroundThreshold 20
  -maxVariation 5
  -maxHeight 150 |
flip -flipCode 0 |
transpose |
houghTextLine
  -resolution 60
  -noTextLineOperation bottom
  -startLambdaBandWidth 0.65
  -endLambdaBandWidth 0.75
  -startLambdaBandMin 0.1
  -endLambdaBandMin 0.0 |
textSkewCorrection
  -maxDegree 52
  -fromDegree 0
  -toDegree 50 |
removeUnderline
  -minWidth 45
  -foregroundThreshold 20
  -maxRelHeight 0.6
  -maxVariation 5 |
flip -flipCode 0 |
removeUnderline
  -minWidth 75
  -foregroundThreshold 20
  -maxRelHeight 0.5
  -maxVariation 5 |
flip -flipCode 0 |
polynomialTextLine
  -operation align
  -order 3
  -outlierFactor 1.0 |
houghTextLine
  -resolution 100

```

```

        -noTextLineOperation original
        -startLambdaBandWidth 0.65
        -endLambdaBandWidth 0.75
        -startLambdaBandMin 0.1
        -endLambdaBandMin 0.0
        -operation align
        -deleteAboveAscenders
        -deleteBelowDescenders |
normalizeText
        -belowBaseline 30
        -aboveBaseline 48
        -minCroppedAboveRatio 0.0
        -minCroppedBelowRatio 0.0 |
normalize

```

This pipeline was used for pixel-based feature extraction:

```

grayscale |
flip |
convertToFloat |
normalize -newMax 1 |
featExtract
        -winWidth 3
        -winShift 2
        -featRawCellHeight 1
        -featRawCellWidth 1
        -featRawCellShift 1
        -kaldiFile data/feats/fsushi.ark,t |
devNull

```

The segment-based features were generated with the following pipeline.

```

grayscale |
flip |
convertToFloat |
normalize -newMax 1 |
featExtract
        -extractors snake
        -winWidth 3
        -winShift 2
        -kaldiFile data/feats/fslytherin.ark,t |
devNull

```

2.4.2 Document Skew Detection Based on Hough Space Derivatives

We presented in Stahlberg and Vogel (2015a) a novel method for document skew estimation using gradients in the Hough transformed image. The `exactOrientationCorrection` operation in `PrepOCReprocessor` is an implementation of the described procedure. The following pipeline has been used in Stahlberg and Vogel (2015a):

```

grayscale |
threshold -type BINARY_INV,OTSU |
extendForHough -maxAngle 15.0 |
exactOrientationCorrection
        -maxAngle 17.0
        -resolution 60

```

```
-eps 0.001
-criterion sum
-noCorrection
```

2.4.3 Detecting Dense Foreground Stripes in Arabic Handwriting for Accurate Baseline Positioning

Stahlberg and Vogel (2015b) describes our approach for detecting the baseline in Arabic handwritings. The following pipeline corresponds to the best result on the IFN/ENIT database reported in this paper:

```
grayscale |
threshold -type BINARY_INV |
removeDiacritics |
houghTextLine
    -startLambdaBandWidth 0.25
    -endLambdaBandWidth 0.45
    -startLambdaBandMin 0.2
    -endLambdaBandMin 0.4
    -combination lowest
    -resolution 130
```

For the KHATT corpus, we have to deal with curved and discontinuous baselines. We use the following pipeline for the KHATT corpus. It also deals with background borders around the text.

```
grayscale |
threshold -type BINARY,OTSU |
vertTextSegmentation
    -minSlope 0.00000001
    -minMargin 0.2
    -concatChildren |
transpose |
vertTextSegmentation
    -minSlope 0.00000001
    -minMargin 0.2
    -concatChildren |
transpose |
houghTextLine |
transpose |
vertTextSegmentation
    -minSlope 0.00000001
    -minMargin 0.2
    -concatChildren |
transpose |
splitTextLines -minWidth 2.5 |
houghTextLine |
concat |
transpose |
vertTextSegmentation
    -minSlope 0.00001
    -minMargin 0.2
    -concatChildren |
transpose |
houghTextLine
```

3 Troubleshooting

This section contains some of the most common error messages and their solutions.

→ Section 2.1

Exception in thread "main" java.lang.NoClassDefFoundError: org.opencv.core.Core This error usually occurs if the path to the OpenCV jar file is not set correctly. Check the `OPENCV_JAR_PATH` variable in the `prepocressor` file in the root directory of your PrepOCReSSor installation. You can find detailed installation instructions in Section 2.1.

→ Section 2.1

Exception in thread "main" java.lang.UnsatisfiedLinkError: no opencv_java2410 in java.library.path This error usually occurs if the path to the OpenCV native library is not set correctly. Check the `OPENCV_NATIVE_LIB` variable in the `prepocressor` file in the root directory of your PrepOCReSSor installation. The variable must point to a directory containing a file called `opencv_java24x.so` where `x` corresponds to your OpenCV version (2.4.x). You can find detailed installation instructions in Section 2.1.

OpenCV Error: Assertion failed (src.type() == CV_XY) `XY` stands for a certain data type (like `8UC1`, `32FC3`, see the OpenCV documentation for more information). This error indicates that two consecutive operations in the pipeline do not fit together. For example, the following command usually results in such a type error.

```
prepocressor -inputFile test.jpg -pipeline "threshold"
```

→ Section 4.2.23

→ Section 4.2.8

The reason is that the input image usually consists of three channels, and the `threshold` operation expects a single channel image. A preceding `grayscale` operation (i.e. the pipeline "`grayscale|threshold`") results in the expected behaviour. Useful operations for resolving this type of error are the `grayscale` operation (Section 4.2.23) and the `convertToFloat` operation (Section 4.2.8).

Memory consumption explodes with large batch sizes Early PrepOCReSSor versions had problems with memory leaks resulting in a linear increase of required RAM in time. However, as long as the images in the batch have approximately the same size/complexity, the memory requirement should be constant after the initial start-up phase. If you have problems with memory, please report the pipeline to us and we try to fix the issue. In the meantime, it helps to split the batch input file into multiple smaller files (e.g. using the Linux `split` command) and call PrepOCReSSor for each split file separately.

```
man split
```

4 Operation and Parameter Reference

4.1 Global Parameters

PREPOCRESSOR 0.2 is a tool for preprocessing images and feature extraction for OCR developed at the Qatar Computing Research Institute. Configuration via command line arguments: `-<name> <value>` Configuration via file (per line): `<name> <value>` Following parameters are available:

<code>configDumpFileName</code>	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
<code>configFile</code>	<i>String, Default: <not set></i> – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
<code>idLength</code>	<i>Integer, Default: 3</i> – Length of the numerical IDs that are inserted when one image in the pipeline produces multiple children. Fill up with trailing 0s if the number is shorter. Set to 0 to switch off trailing 0s.
<code>inputFile</code>	<i>String, Default: imageList.txt</i> – Text file containing paths to the input images. The paths should be separated by line breaks. This parameter can also point directly to an image file if only one image is to be processed. Following input formats are supported (provided by OpenCV's <code>imread</code> function): <ul style="list-style-type: none"> • Windows bitmaps - *.bmp, *.dib • JPEG files - *.jpeg, *.jpg, *.jpe • JPEG 2000 files - *.jp2 • Portable Network Graphics - *.png • Portable image format - *.pbm, *.pgm, *.ppm • Sun rasters - *.sr, *.ras • TIFF files - *.tiff, *.tif • Additionally, the CSV file format is supported by preprocessor.
<code>logLevel</code>	<i>String, Default: INFO</i> – Controls the amount of output. <ul style="list-style-type: none"> • FATAL: Only fatal errors, • ERROR: All errors, • WARN: Warnings and errors, • INFO: Notices, warnings and errors, • DEBUG: Debug mode
<code>nThreads</code>	<i>Integer, Default: 1</i> – Number of threads.
<code>outputPath</code>	<i>String, Default: 15-06-24.105233.felix-%base%idx%ext</i> – This parameter controls where the output files are stored. The file format is determined by the file extension defined by this template. For available file formats, see the <code>-inputFile</code> parameter. Following placeholders can be used: <ul style="list-style-type: none"> • '%dir': Name of the directory containing the input image. • '%base': Base name of the input image file name (without file extension). • '%-base': Same as %base, but cut after first minus hyphn ('-') • '%idx': Some operators split up input images into smaller pieces. The pieces are stored substituting %idx with '1', '2'... • '%unqatip': Assume base name of the input image file name (without file extension) is a QATIP style ID. Then the %unqatip placeholder stands for the original corpus id. Breaks if original speaker or utterance id starts with 'x'

- '%ext': File extension (given by the input file).

pipeline *String, Default: <not set>* – This parameter defines the operations to be executed on the input images. The syntax is similar to the linux shell pipeline: Operations are separated by '—' und parameterized with the common-<arg> <val> syntax. For details regarding the operations, try -help <operation-name>. Available operations are:

- adaptiveThreshold
- axisAlignedHough
- bbq
- blur
- col2graph
- componentDensity
- concat
- convertToFloat
- cutWithAltecXml
- cutWithPageXml
- devNull
- drawChildren
- drawKaldiAlignment
- drawTextLines
- exactOrientationCorrection
- extractConstantRegions
- extend
- extendForHoughsquare
- featExtract
- fillTransparency
- filter
- flip
- grayscale
- hough
- houghTextLine
- invert
- log
- morph
- multiChannelOtsu
- normalize
- normalizeText
- normalizeUpperBaseline
- orientationCorrection

- outlierRemove
- polynomialTextLine
- printMax
- projectionLineSegmentation
- rectSum
- reduce
- reducedAlcmTransform
- renderPageXmlTranscriptions
- removeDiacritics
- removeLargeComponents
- removeSmallComponents
- removeUnderline
- removeVertTextMargin
- scale
- sobel
- splitTextLines
- subtractMean
- tee
- textSkewCorrection
- thinning
- threshold
- transpose
- vertTextSegmentation
- writeRects

- `silentOverwrite` *Integer, Default: 1* – No files are overridden if this is set to 0.
- `singlePopulation` *Integer, Default: 0* – If this parameter is set to 1, only a single population is used for all input images. Otherwise, a population is created for each of the input files separately. If this option is set the `-nThreads` parameter is not used. Note that the file name of the population is set to the first input file

4.2 Operations

4.2.1 adaptiveThreshold Operation

The `adaptiveThreshold` command creates binary images. It is based on the OpenCV function `adaptiveThreshold()`. Following parameters are available:

- C** *Float, Default: 2.0* – C constant (passed through to OpenCV).
- `adaptiveType` *String, Default: MEAN_C* – Adaptive thresholding method. See OpenCV's documentation for the `adaptiveThreshold` function. Available values are `MEAN_C` or `GAUSSIAN_C`

<code>blockSize</code>	<i>Integer, Default: 12</i> – Block size (passed through to OpenCV).
<code>configDumpFileName</code>	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
<code>configFile</code>	<i>String, Default: <not set></i> – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
<code>maxVal</code>	<i>Integer, Default: 255</i> – Lowest possible value (passed through to OpenCV).
<code>type</code>	<i>String, Default: BINARY</i> – Thresholding type. See OpenCV's documentation for the threshold function. Connect options with ','. Available options are: BINARY, BINARY_INV, TRUNC, TOZERO, TOZERO_INV, OTSU

4.2.2 axisAlignedHough Operation

This is a specialized and modified version of the Hough transformation. In contrast to the Hough space, rho is always on the x-axis. Rho ranges from 0 to image width. The range and resolution for theta can be specified. The advantage of this implementation is that there are no quantization errors for rho since the resolution is exactly one pixel. The disadvantage is that only lines crossing the x axis between 0 and image width are considered. The returned image contains the counts where the y coordinate represents theta. Following parameters are available:

<code>configDumpFileName</code>	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
<code>configFile</code>	<i>String, Default: <not set></i> – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
<code>fromTheta</code>	<i>Float, Default: 45.0</i> – Minimum value for theta.
<code>thetaResolution</code>	<i>Integer, Default: 90</i> – Number of theta quantization steps.
<code>toTheta</code>	<i>Float, Default: 135.0</i> – Maximum value for theta.

4.2.3 bbq Operation

Keep only the lowest point in each column. The lowest point is detected by comparing the first channel with the threshold parameter. Following parameters are available:

<code>configDumpFileName</code>	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
<code>configFile</code>	<i>String, Default: <not set></i> – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
<code>threshold</code>	<i>Float, Default: 0.5</i> – Threshold for detecting the lowest point

4.2.4 blur Operation

Blurs the images. Following parameters are available:

- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- mode** *String, Default: mean* – Blur mode. Available: mean, gaussian, median. For the median filter, -xSize is used for both dimensions.
- xSize** *Integer, Default: 5* – Kernel size in x direction.
- ySize** *Integer, Default: 5* – Kernel size in y direction.

4.2.5 col2graph Operation

Converts the first column of the image to a graph. Following parameters are available:

- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- graphWidth** *Integer, Default: 100* – Width of the generated image.

4.2.6 componentDensity Operation

Calculates a map of connected component density. Each connected component adds 1/area to all pixels in its bounding box where area is the size of the bounding box. Following parameters are available:

- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- minSize** *Integer, Default: 4* – Bounding boxes below this value are ignored.

4.2.7 concat Operation

Concatenate all images of one population horizontally. Following parameters are available:

- center** *Integer, Default: 1* – Set to 0 if the images should not be centered in case of different heights
- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file.

Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

4.2.8 convertToFloat Operation

Convert Matrix to CV_32FC1. Following parameters are available:

- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

4.2.9 cutWithAltecXml Operation

This operation corresponds to cutWithPageXml but reads xml files in the format used by the ALTEC corpus. This format specifies line tags on the first level and word tags on the second level. See <http://ALTEC-Center.org/xsd/ocr-annotation-1-0.xsd> for a specification. Following parameters are available:

- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- cutLevel** *String, Default: line* – Splitting level. Either 'line' or 'word'
- useIndexAttributes** *Integer, Default: 0* – Set the preprocessor index to the value of the index attribute of the corresponding node in the xml file. Note that this can lead to problems in combination with cutLevel=word because the same word level index might be used multiple times in a single xml file.
- xmlPath** *String, Default: %base%idx.xml* – Path to the xml files in ALTEC format. The same placeholders as in the global outputPath can be used.

4.2.10 cutWithPageXml Operation

Cut a page image using an XML file in PAGE format. Note: The used PAGE library may break with multiple threads! Following parameters are available:

- border** *Integer, Default: 0* – If -border equals 0 we cut the text regions accurately. Use a value greater than zero if the extracted regions are used for human inspections. It will add a padding to the image with decreased brightness and draw a red rectangle around the region.
- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- extractRegions** *Integer, Default: 1* – Extract regions.

- extractTextObjects** *Integer, Default: 0* – Extract text lines.
- minLevel** *Integer, Default: 0* – Minimum level in layout of region to be extracted.
- usePageIds** *Integer, Default: 0* – Use id attributes in xml file for naming. Otherwise, use consecutive numbering (see global idLength parameter)
- xmlPath** *String, Default: %base%idx.xml* – Path to the xml files in PAGE format. The same placeholders as in the global outputPath can be used.

4.2.11 devNull Operation

Deletes all images in the pipeline. Equivalent to '> /dev/null' in the unix shell. Following parameters are available:

- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

4.2.12 drawChildren Operation

Reloads the input image of the given population and draws all children in the population with rectangles. Following parameters are available:

- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- thickness** *Integer, Default: 2* – Thickness of the rectangle. Negative for filled rectangles.
- transpose** *Integer, Default: 0* – Transpose original.

4.2.13 drawKaldiAlignment Operation

Reads a Kaldi alignment file and draws the forced alignment into the image. Assumes left-to-right topology for nonsilence phones and whatever for silence. Following parameters are available:

- alignmentFile** *String, Default: kaldi.al* – Path to the Kaldi alignment file in text format
- borderHeight** *Integer, Default: 30* – Height of the border for annotations.
- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- kaldiId** *String, Default: %base* – This string specifies how the kaldi ID is generated. You can use the same placeholders as in outputPath.
- offset** *Integer, Default: 0* – Offset from the right image border.

4.2.14 drawTextLines Operation

Reloads the original input images and draws text lines in it for visual verification. The text lines must be in the pipeline, e.g. produced by the houghTextLine operation. Following parameters are available:

<code>configDumpFileName</code>	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
<code>configFile</code>	<i>String, Default: <not set></i> – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
<code>drawSkewLine</code>	<i>Integer, Default: 1</i> – Draw the skew line. Set to 0 to ignore the text skew information

4.2.15 exactOrientationCorrection Operation

Brings rotated text documents in an upright position. This is done by finding the maximum squared variance angle in the Hough transformed image. An iterative algorithm is applied to increase the accuracy of the skew angle estimation. Following parameters are available:

<code>configDumpFileName</code>	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
<code>configFile</code>	<i>String, Default: <not set></i> – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
<code>criterion</code>	<i>String, Default: horiz</i> – Maximization criterion. Available values: <ul style="list-style-type: none">• 'horiz': Horizontal estimation• 'vert': Vertical estimation• 'sum': Sum of horizontal and vertical estimation
<code>eps</code>	<i>Float, Default: 0.1</i> – Accuracy in degree.
<code>horizWeight</code>	<i>Float, Default: 0.5</i> – If the sum criterion is used, the horizontal profile is weighted with horizWeight and the vertical profile is weighted with (1-horizWeight)
<code>houghLineMode</code>	<i>String, Default: scaling</i> – Method for line definition in Hough space. Available values: <ul style="list-style-type: none">• 'scaling': Gradually increase scaling factor of line definition• 'bresenham': Use Bresenham's line drawing algorithm• 'exact': Take fractional counts for pixels into account
<code>maxAngle</code>	<i>Float, Default: 45.0</i> – Maximum skew angle.
<code>noCorrection</code>	<i>Integer, Default: 0</i> – Pass thru image without modification.
<code>refine</code>	<i>Integer, Default: 0</i> – Set to 1 to enable refinement. If this parameter is set an additional search in degree +/- 0.5 with resolution 100 is added assuring that the tested values are multipliers of 0.01
<code>reloadOriginal</code>	<i>Integer, Default: 0</i> – Reload the original image and rotate it. Otherwise use image in the pipeline.

- resolution** *Integer, Default: 90* – Resolution of the Hough transform.
- sobelKSize** *Integer, Default: 3* – Size of the sobel kernel.

4.2.16 extractConstantRegions Operation

This operations assumes that the images in the pipeline are frames of a video. Note that the fps video filter in ffmpeg can be used to extract an image every X seconds of the video. The operation searches for regions which do not change during a certain time period. In news videos, these regions usually correspond to overlays embedded in the video displaying additional information. In the documentation for this operation we refer a single image in the pipeline as 'frame'. Following parameters are available:

- closeSize** *Integer, Default: 30* – Kernel size of closing operation
- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- minLength** *Integer, Default: 5* – Minimal time period for a region in number of frames
- openSize** *Integer, Default: 150* – Kernel size of open operation
- threshold** *Float, Default: 20.0* – Maximal sum of differences in channels for a pixel to be considered as constant.

4.2.17 extend Operation

Extend image canvas. Following parameters are available:

- bottom** *Integer, Default: 20* – Extend image at bottom border (in pixels).
- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- left** *Integer, Default: 20* – Extend image at left border (in pixels).
- right** *Integer, Default: 20* – Extend image at right border (in pixels).
- top** *Integer, Default: 20* – Extend image at top border (in pixels).

4.2.18 extendForHoughsquare Operation

10:52:51 FATAL: Operation 'extendForHoughsquare' is not implemented

4.2.19 featExtract Operation

Feature extraction for Kaldi. The feat* parameters are extractor specific. NOTE: Feature extraction is based on a sliding window in right-to-left direction as this tool was initially developed for Arabic. If you wish to

change direction, apply the flip operation first. Following parameters are available:

- baselineHeight** *Integer, Default: 32* – Height of the baseline for baseline dependent features.
- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- delayDelta** *Integer, Default: 1* – Set to positive value to add deltas of feature vectors according -delays
- delayRaw** *Integer, Default: 0* – Set to positive value to add raw feature vectors according -delays (feature staking)
- delays** *String, Default: <not set>* – Comma-separated list of integers specifying the deltas to add. The integers are the delta distances, i.e. '1' stands for standard deltas, '2' calculates deltas to second last feature vector.
- extractors** *String, Default: raw* – Comma separated list of feature extractors. Available:
- 'raw': Raw pixel values.
 - 'directional': Directional features.
 - 'snake': Snake feature extraction method (also called segment-based method in (Stahlberg and Vogel, 2015))
 - 'runlengths': Use pixel-wise runlengths in 4 directions
 - 'anhdf': ANHDF features (El-Mahallawy, 2008)
 - 'distribution': Distribution features as defined by (Likforman-Sulem et. al., 2012)
 - 'concavity': Concavity features as defined by (Likforman-Sulem et. al., 2012)
- featAnhdfConnectivityTolerance** *Integer, Default: 4* – Tolerance for segments in the ANHDF feature to be connected. See (El-Mahallawy, 2008) PhD thesis for more information.
- featAnhdfReductionMode** *String, Default: max* – ANHDF features are defined for windows with 1 pixel width. Wider windows a reduced according to this method:
- 'max': Take the maximum of each row.
 - 'min': Take the minimum of each row.
 - 'average': Take the average of each row.
 - 'firstAndLast': Use right most column for slice i and left most column of previous slice as i-1.
- featAnhdfSegmentNum** *Integer, Default: 4* – Number of segments in ANHDF features. 4 is also used by (El-Mahallawy, 2008) and reasonable for Arabic.
- featConcavityBaselineDependent** *Integer, Default: 1* – Extract also concavity separately for above and below baseline.
- featDirectionalRadius** *Integer, Default: 10* – Radius for directional feature extractor (maximum feature value)

<code>featRawCellHeight</code>	<i>Integer, Default: 1</i> – Height of the cell for the raw feature extractor.
<code>featRawCellShift</code>	<i>Integer, Default: 1</i> – Vertical cell shift for the raw feature extractor.
<code>featRawCellWidth</code>	<i>Integer, Default: 1</i> – Height of the cell for the raw feature extractor.
<code>featRunlengthsNonNegative</code>	<i>Integer, Default: 1</i> – Set to positive value to avoid using negative values for background pixel runlengths (use 0 instead)
<code>featRunlengthsRadius</code>	<i>Integer, Default: 10</i> – Radius for runlength feature extractor (maximum feature value)
<code>featSnakeAddCenterDistances</code>	<i>Integer, Default: 0</i> – Add distance between consecutive segment centers as features.
<code>featSnakeAddRelativeFeats</code>	<i>Integer, Default: 0</i> – Add snake features divided by height of entire slice
<code>featSnakeBackground</code>	<i>Integer, Default: 0</i> – Set to positive value to use background snakes.
<code>featSnakeDefaultHeight</code>	<i>Integer, Default: 0</i> – Default value for height features which is used in silence.
<code>featSnakeForeground</code>	<i>Integer, Default: 1</i> – Set to positive value to use foreground snakes.
<code>featSnakeNumber</code>	<i>Integer, Default: 6</i> – Number of snakes.
<code>foregroundThreshold</code>	<i>Float, Default: 0.5</i> – Pixel values above this threshold are considered as foreground.
<code>kaldiFile</code>	<i>String, Default: kaldi.ark,t</i> – Path to the kaldi feature file to generate. This is a feature table in text format. See kaldi's copy-feat tool with ark,t specifiers for more information.
<code>kaldiId</code>	<i>String, Default: %base</i> – This string specifies how the kaldi ID is generated. You can use the same placeholders as in <code>outputPath</code> .
<code>winShift</code>	<i>Integer, Default: 2</i> – Shift of the sliding window.
<code>winWidth</code>	<i>Integer, Default: 3</i> – Width of the sliding window.

4.2.20 fillTransparency Operation

Loads the original image, fetches the alpha channel and fills transparent areas in the current image in the pipeline with zero. Following parameters are available:

<code>configDumpFileName</code>	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
<code>configFile</code>	<i>String, Default: <not set></i> – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

4.2.21 filter Operation

Removes images that are likely to be no text lines. Following parameters are available:

<code>configDumpFileName</code>	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
<code>configFile</code>	<i>String, Default: <not set></i> – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

- minAspectRatio** *Float, Default: 2.0* – Images with width/height<minAspectRatio are removed.
- minHeight** *Integer, Default: 10* – Images smaller height (in pixel) are removed.

4.2.22 flip Operation

Flips the image around x or y axis. Following parameters are available:

- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- flipCode** *Integer, Default: 1* – Except from OpenCV docu: Specifies how to flip the array: 0 means flipping around the x-axis, positive (e.g., 1) means flipping around y-axis, and negative (e.g., -1) means flipping around both axes.

4.2.23 grayscale Operation

Converts the images in the pipeline to grayscale. Following parameters are available:

- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

4.2.24 hough Operation

The hough command performs a Hough transformation on the input images. This is the original OpenCV implementation. Consider using axisAlignedHough if this is not the last operation in the pipeline or you want to obtain a meaningful image. This operation is useful for directly redirecting the output to a CSV file. Following parameters are available:

- angleResolution** *Integer, Default: 360* – Angle resolution (number of different values for theta)
- angleSamplingFactor** *Integer, Default: 200* – HoughLines is called with resolution 1/angleSamplingFactor. A higher value reduces the noise in the Hough transform, but needs longer execution time.
- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- toMatrix** *Integer, Default: 0* – Transform the hough transformation to a matrix. The rho values are shifted so that the minimum value corresponds to x=0, and x=width/2 to rho=0. The theta values are scaled by angleResolution.

If toMatrix is not set, the operation passes an array of 3 dimensional vectors storing [rho, theta, voteCount] ordered descending by voteCount.

4.2.25 houghTextLine Operation

Calculates the base line from a Hough transformed image. If -criterion=max, the base line is detected at the maximum in the Hough space. Otherwise, we find a rotated rectangle which includes at least minBetweenBaseline white pixels and optimizes the target function (-criterion parameter) and meets the -maxValArea restriction. The base line is detected at the bottom of the rectangle. See (Stahlberg and Vogel, 2015) for a detailed discussion. Following parameters are available:

blurMode	<i>String, Default: none</i> – Blur hough space. Available values are: 'none', 'median', 'mean'.
blurRho	<i>Integer, Default: 3</i> – Kernel size of blur operation in rho direction
blurTheta	<i>Integer, Default: 1</i> – Kernel size of blur operation in theta direction
borderFactor	<i>Float, Default: 0.0</i> – Before Hough transform, the image is extended on top and bottom by borderFactor*height
combination	<i>String, Default: lowest</i> – How the different baselines should be combined. Available: 'none': Do not combine baselines - Pass them separately 'lowest': Select the baseline with the lowest slope.
configDumpFileName	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
configFile	<i>String, Default: <not set></i> – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
deleteAboveAscenders	<i>Integer, Default: 0</i> – Set to 1 to remove pixels above the highest ascender. An ascender is a connected component which crosses the upper baseline. Ignored if the operation parameter is not set to 'align'
deleteBelowDescenders	<i>Integer, Default: 0</i> – Set to 1 to remove pixels below the lowest descender. A descender is a connected component which crosses the lower baseline. Ignored if the operation parameter is not set to 'align'
deltaLambda	<i>Float, Default: 0.05</i> – Lambda increment. Set to 0 to use only start-Lambda*.
endLambdaBandMin	<i>Float, Default: 0.4</i> – Required fraction of white pixels between both base lines for the bandMin criterion.
endLambdaBandWidth	<i>Float, Default: 0.4</i> – Required fraction of white pixels between both base lines for the bandWidth criterion.
houghMax	<i>Integer, Default: 0</i> – Set to 1 to include the hough space maximum.
maxDegree	<i>Float, Default: 20.0</i> – Maximum text rotation in degree. If this rotation is exceeded, the image is discarded. Note: We only check the range between +-45 degree.
maxValArea	<i>Float, Default: 0.5</i> – The maximum within the band must be in the lower part of the band area. This restriction addresses the common assumption that the baseline is represented by a maximum in the Hough space. Set to negative value to disable this check. This is not used if criterion=max
noTextLineOperation	<i>String, Default: original</i> – This parameter decides what is passed through

the pipeline if no text line within range has been detected:

- 'original': Pass through original image.
- 'bottom': Place baseline at bottom border of image

operation	<i>String, Default: align</i> – This parameter decides what is passed through the pipeline. Available values are: <ul style="list-style-type: none">• 'none': Leave the images as they are.• 'draw': Draw upper and lower baselines.• 'align': Create a new image where the baseline is horizontal and at image height/2
resolution	<i>Integer, Default: 130</i> – Number of steps between -45 and +45 degree in Hough transform.
startLambdaBandMin	<i>Float, Default: 0.2</i> – Required fraction of white pixels between both base lines for the bandMin criterion.
startLambdaBandWidth	<i>Float, Default: 0.2</i> – Required fraction of white pixels between both base lines for the bandWidth criterion.
truPath	<i>String, Default: <not set></i> – If .tru files are available (as for the IFN/ENIT database) the baseline error can be computed automatically. Use placeholders as in the global -outputPath parameter. The evaluation is written to stdout. Format: EvalBaseline <inputFileName> <refStartY> <refSlope> <hypoStartY> <hypoSlope> <Error> <StringError>
useTruIfAvailable	<i>Integer, Default: 1</i> – If this is set to 1, we take the reference baseline from the tru file if exists.

4.2.26 invert Operation

Invert the image. Following parameters are available:

configDumpFileName	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
configFile	<i>String, Default: <not set></i> – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

4.2.27 log Operation

Element wise logarithm. Following parameters are available:

configDumpFileName	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
configFile	<i>String, Default: <not set></i> – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

4.2.28 morph Operation

Performs morphological operations. Following parameters are available:

configDumpFileName	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
---------------------------	---

- `configFile` *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- `kernelShape` *String, Default: rect* – Kernel shape. 'rect' or 'ellipse'
- `kernelSize` *Integer, Default: 5* – Kernel size
- `operation` *String, Default: close* – One of the following morphology operations: close, open, erode, dilate

4.2.29 multiChannelOtsu Operation

This is Otsu thresholding adapted for multichannel images. It uses greyscale standard otsu binarization for initial labeling, and then applies the k-means algorithm (k=2) for final binarization. Note: Channels > 3 (e.g. alpha channel) are not considered. Following parameters are available:

- `blackDiscount` *Float, Default: 0.5* – Increase this to make more pixels classified as white. Between 0 and 1
- `configDumpFileName` *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- `configFile` *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- `extractRegions` *Integer, Default: 0* – Extract regions.
- `extractTextObjects` *Integer, Default: 1* – Extract text lines.
- `maxForegroundFraction` *Float, Default: 0.2* – Maximum fraction of foreground pixels. If exceeded, increase blackDiscount parameter to produce more background
- `maxIter` *Integer, Default: 10* – Number of k-means iterations
- `minLevel` *Integer, Default: 0* – Minimum level in layout of region to be extracted.
- `normalizeRegionChannels` *Integer, Default: 0* – Normalize channels in top level regions before binarization.
- `xmlPath` *String, Default: <not set>* – Path to the xml files in PAGE format. The same placeholders as in the global outputPath can be used. If this parameter is provided, binarization is done for each region separately. Otherwise, the algorithm is applied to the whole image. Note: The used PAGE library may break when using multiple threads

4.2.30 normalize Operation

The normalize operation rescales the values in the matrices to the given interval. The current value range is fetched from the first channel only. Following parameters are available:

- `configDumpFileName` *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- `configFile` *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

newMax *Float, Default: 255.0* – Maximum value of the new interval.

newMin *Float, Default: 0.0* – Minimum value of the new interval.

4.2.31 `normalizeText` Operation

Assumes that images contain text lines with horizontal baseline at image center. Scales images such that the baseline is repositioned as defined by `-below/aboveBaseline` and scaled such that the first/last row is the nearest row to the baseline which sums up to less than `-maxCut` pixels. The resulting images will have the height `belowBaseline+aboveBaseline`. Following parameters are available:

aboveBaseline *Integer, Default: 32* – Vertical distance in the output image from baseline to upper border.

belowBaseline *Integer, Default: 16* – Vertical distance in the output image from baseline to lower border.

configDumpFileName *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.

configFile *String, Default: <not set>* – Configuration file (format: `<key> <val>`). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

maxBelowShrink *Float, Default: 3.0* – If scale below and above baseline differ by this factor (relative) use above bl scale for below bl.

maxBelowStretch *Float, Default: 1.5* – If scale below and above baseline differ by this factor (relative) use above bl scale for below bl.

maxCut *Float, Default: 1.0* – Maximum sum at cropped border

minCroppedAboveAbsolute *Integer, Default: 20* – Distance of cropped border above baseline.

minCroppedAboveRatio *Float, Default: 0.1* – Distance of cropped border above baseline (relative to image border).

minCroppedBelowRatio *Float, Default: 0.02* – Distance of cropped border above baseline (relative to image border).

4.2.32 `normalizeUpperBaseline` Operation

This operation normalizes the position of the upper baseline. The input image should be an aligned image with lower baseline in image center (see `houghTextLine` operation) The upper baseline is estimated at the maximum in the derivative of the horizontal projection profile above the lower baseline. The image is modified s.t. the upper baseline is at a predefined height. Note: If you apply the `normalizeText` operation after this, the `maxCut`, `minCroppedAboveRatio`, and `minCroppedAboveAbsolute` parameters should be equal. Following parameters are available:

configDumpFileName *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.

configFile *String, Default: <not set>* – Configuration file (format: `<key> <val>`). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

keepCoreZoneAspectRatio *Integer, Default: 0* – Set to 1 to keep the aspect ratio in the core

zone between upper and lower baseline. Otherwise, the core zone is stretched/shrunked in order to reposition the upper baseline. Only applicable if operation=align

<code>maxCut</code>	<i>Float, Default: 1.0</i> – Maximum sum at cropped border
<code>maxStretchFactor</code>	<i>Float, Default: 4.0</i> – Works with keepCoreZoneAspectRatio=1. Maximum horizontal stretching factor
<code>minCroppedAboveAbsolute</code>	<i>Integer, Default: 20</i> – Distance of cropped border above baseline.
<code>minCroppedAboveRatio</code>	<i>Float, Default: 0.5</i> – Distance of cropped border above baseline (relative to image border).
<code>minStretchFactor</code>	<i>Float, Default: 0.25</i> – Works with keepCoreZoneAspectRatio=1. Minimum horizontal stretching factor
<code>newUpperBaseline</code>	<i>Float, Default: 0.4</i> – New ratio between distance between baselines and highest ascender - lower baseline distance. This is only applicable in combination with -operation=align
<code>operation</code>	<i>String, Default: align</i> – What should be done after the upper baseline is found 'align': Reposition the baseline to a predefined height 'draw': Draw a line indicating the upper baseline
<code>upperBaselineHighest</code>	<i>Float, Default: 0.8</i> – Highest possible ratio between distance between baselines and highest ascender - lower baseline distance
<code>upperBaselineLowest</code>	<i>Float, Default: 0.2</i> – Lowest possible ratio between distance between baselines and highest ascender - lower baseline distance

4.2.33 orientationCorrection Operation

Brings rotated text documents in an upright position. This is done by finding the maximum squared variance angle in the Hough transformed image. Following parameters are available:

<code>configDumpFileName</code>	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
<code>configFile</code>	<i>String, Default: <not set></i> – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
<code>maxAngle</code>	<i>Float, Default: 45.0</i> – Maximum skew angle.
<code>noCorrection</code>	<i>Integer, Default: 0</i> – Pass thru image without modification.
<code>reloadOriginal</code>	<i>Integer, Default: 0</i> – Reload the original image and rotate it. Otherwise use image in the pipeline.
<code>resolution</code>	<i>Integer, Default: 90</i> – Resolution of the Hough transform.
<code>sobelKSize</code>	<i>Integer, Default: 3</i> – Size of the sobel kernel.

4.2.34 outlierRemove Operation

Removes outlier. Outlier are identified by differing by -tolerance times standard derivation from the mean. Input needs to be 1 channel float. Following parameters are available:

<code>configDumpFileName</code>	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
---------------------------------	---

- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- tolerance** *Float, Default: 3.0* – Tolerance parameter for outlier detection.

4.2.35 polynomialTextLine Operation

Fits a polynomial to the data points in order to guess the text line. Following parameters are available:

- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- dataPoints** *String, Default: minima* – Strategy for data point retrieval. 'bbq': Use all bottom foreground points (see bbq op) 'minima': Use only minima of bottom foreground points
- minHeight** *Integer, Default: 1* – Minimum distance to top border for a polynomial in order to be used for estimating the polynomial
- operation** *String, Default: align* – This parameter decides what is passed through the pipeline. Available values are:
- 'none': Leave the images as they are.
 - 'draw': Draw upper and lower baselines.
 - 'align': Create a new image where the baseline is horizontal and at image height/2
- order** *Integer, Default: 4* – Order of the polynomial
- outlierFactor** *Float, Default: -1.0* – Remove data points this factor times stdDev from median. Set to negative value to disable outlier detection
- threshold** *Float, Default: 0.5* – Threshold for detecting the lowest point

4.2.36 printMax Operation

Print information about the maximum. Following parameters are available:

- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

4.2.37 projectionLineSegmentation Operation

Line segmentation using vertical projection. This operation first fits a sinus function to the profile. The frequency of the best fit is then used

to determine the kernel size for a blur operation on the projection profile. Lines are extracted from valley to valley in the smoothed profile. Following parameters are available:

<code>analysisMode</code>	<i>Integer, Default: 0</i> – Set to 1 to output an image explaining the line segmentation by showing the smoothed vertical projection over the image plus the found boundaries
<code>configDumpFileName</code>	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
<code>configFile</code>	<i>String, Default: <not set></i> – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
<code>maxLineHeight</code>	<i>Integer, Default: 130</i> – Maximum line height in pixel.
<code>maxLineHeightVariance</code>	<i>Float, Default: 2.0</i> – If a line is taller than this parameter times the estimated average line height -> outlier.
<code>maxProjectionRatio</code>	<i>Float, Default: 0.8</i> – If the minimum next to a maximum is larger than this parameter times the maximum in the projection, ignore this maximum.
<code>minLineCount</code>	<i>Integer, Default: 10000</i> – Minimal line count in one segment. This can be used for outlier detection. Set to a high value to disable this feature.
<code>minLineHeight</code>	<i>Integer, Default: 10</i> – Minimal line height in pixel.
<code>sinusExp</code>	<i>Integer, Default: 1</i> – Sinus exponent used for line height estimation. Should be uneven.

4.2.38 rectSum Operation

This operation calculates the sum of elements within rectangles in the images. The rectangles have a common edge point (fix point) but their width and height vary. Produces a `maxWidth` times `maxHeight` matrix storing the sum within corresponding rectangles. Following parameters are available:

<code>configDumpFileName</code>	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
<code>configFile</code>	<i>String, Default: <not set></i> – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
<code>fixX</code>	<i>Integer, Default: 0</i> – x coordinate of fix point
<code>fixY</code>	<i>Integer, Default: 0</i> – y coordinate of fix point
<code>maxHeight</code>	<i>Integer, Default: 0</i> – Largest rectangle height. Can also be negative. If it is set to 0, use <code>imageHeight-fixY</code>
<code>maxWidth</code>	<i>Integer, Default: 0</i> – Largest rectangle width. Can also be negative. If it is set to 0, use <code>imageWidth-fixX</code>

4.2.39 reduce Operation

Calculates the projections of the images in the pipeline. The images are reduced to a single column or row (see `dim` parameter). Following parameters are available:

configDumpFileName *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.

configFile *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

dim *Integer, Default: 1* – Dimension along which the reduction is done. E.g. 1 reduces the image to a single column.

mode *String, Default: sum* – Accumulation mode. Available modes: sum, avg, max, min, sqrSum

4.2.40 reducedAlcmTransform Operation

Applies a steerable ellipsoid filter to create an adaptive local connectivity map. The ALCM of each direction is reduced horizontally to a single col. The *i*-th col of the resulting image corresponds to the direction *i*. *i* encodes the angle $180*i/resolution$. If resolution=2, *i*=0 is the horizontal, *i*=1 the vertical ALCM. The resulting image has the width -resolution. Following parameters are available:

configDumpFileName *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.

configFile *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

kHeight *Integer, Default: 6* – Height of the ellipse kernel.

kWidth *Integer, Default: 30* – Width of the ellipse kernel.

resolution *Integer, Default: 2* – See operation description.

4.2.41 renderPageXmlTranscriptions Operation

This operation scans a page xml file for text regions. The text regions are written to the document images in the pipeline by inserting solid rectangles with the text of the xml files in it. For example, this can be used to generate a translated version of the document image after the text has been ocred and translated. Following parameters are available:

align *String, Default: center* – Text alignment, 'right', 'left', or 'center'

bgColorEstimateBorder *Integer, Default: 2* – Controls the way the background color for text areas is estimated. The color is the average of the pixel colors at the border of the text area. This is the thickness of that border.

configDumpFileName *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.

configFile *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

fontFamily *String, Default: Arial* – Font Family. Only used if renderWithOpenCV is not set

maxFontSize *Integer, Default: 100* – Maximum font size. Only used if renderWith-

	OpenCV is not set
<code>minFontSize</code>	<i>Integer, Default: 12</i> – Minimum font size. Only used if <code>renderWithOpenCV</code> is not set
<code>minLevel</code>	<i>Integer, Default: 0</i> – Minimum level in layout of region to be extracted.
<code>renderWithOpenCV</code>	<i>Integer, Default: 0</i> – Set to 1 to use OpenCV's <code>putText</code> function for rendering the text. This is faster but does only support ASCII characters.
<code>xmlPath</code>	<i>String, Default: %base%idx.xml</i> – Path to the xml files in PAGE format. The same placeholders as in the global <code>outputPath</code> can be used.

4.2.42 `removeDiacritics` Operation

Remove diacritics (small and quadractic connected components) in the image Following parameters are available:

<code>configDumpFileName</code>	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
<code>configFile</code>	<i>String, Default: <not set></i> – Configuration file (format: <code><key> <val></code>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
<code>maxHeight</code>	<i>Float, Default: 0.2</i> – Maximum height of a diacritic relative to image height
<code>maxWidth</code>	<i>Float, Default: 0.3</i> – Maximum width of a diacritic relative to image height

4.2.43 `removeLargeComponents` Operation

Remove large connected components in the image – i.e. components which exceed either the maximum width or the maximum height. Following parameters are available:

<code>configDumpFileName</code>	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
<code>configFile</code>	<i>String, Default: <not set></i> – Configuration file (format: <code><key> <val></code>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
<code>maxHeight</code>	<i>Float, Default: 0.3</i> – Maximum height of a connected component relative to the image height.
<code>maxWidth</code>	<i>Float, Default: 0.2</i> – Maximum width of a connected component relative to the image width.

4.2.44 `removeSmallComponents` Operation

Remove small connected components in the image – i.e. components which are smaller than both the minimum width and height. Following parameters are available:

<code>configDumpFileName</code>	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
<code>configFile</code>	<i>String, Default: <not set></i> – Configuration file (format: <code><key> <val></code>). Parameters in this file override values in the default configuration file.

Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

minHeight *Float, Default: 0.3* – Minimum height of a connected component relative to the image height.

minWidth *Float, Default: 0.2* – Minimum width of a connected component relative to the image width.

4.2.45 removeUnderline Operation

Remove underlines in text line images based on bottom point analysis: Record lowest foreground point, look for straight lines, estimate line thickness with median height from segments just above a straight line, override with black. Following parameters are available:

configDumpFileName *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.

configFile *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

foregroundThreshold *Float, Default: 0.5* – Threshold for foreground

maxHeight *Integer, Default: 10000* – Minimum distance from underline to image bottom in pixel

maxRelHeight *Float, Default: 0.6* – Minimum height of the underline relative to the image height

maxThickness *Integer, Default: 6* – Maximum thickness of underline in pixels.

maxVariation *Integer, Default: 6* – Minimum width of connected underline in pixels.

minRelWidth *Float, Default: 0.0* – Minimum width of connected underline relative to image width.

minWidth *Integer, Default: 20* – Minimum width of connected underline in pixels.

thicknessFactor *Float, Default: 1.5* – Underlines are removed up to a thickness of thicknessFactor*median thickness of this segment.

4.2.46 removeVertTextMargin Operation

Removes black space on top and bottom of the image, assuming that it contains only one single text line. Following parameters are available:

configDumpFileName *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.

configFile *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

minRowSum *Float, Default: 3.0* – Rows with less than minPixelCount white pixels are marked as black.

minTextHeight *Integer, Default: 15* – Minimal text height in pixel.

4.2.47 scale Operation

Scales the images in the pipeline. Following parameters are available:

- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- xScale** *Float, Default: 1.0* – Scale in x direction.
- yScale** *Float, Default: 1.0* – Scale in y direction.

4.2.48 sobel Operation

Applies a sobel filter to the image and calculates derivatives in x or y direction. Following parameters are available:

- borderType** *String, Default: constant* – Border type. Available values: constant, replicate
- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- kSize** *Integer, Default: 3* – Size of the sobel kernel.
- xOrder** *Integer, Default: 0* – Order of derivative in x direction.
- yOrder** *Integer, Default: 1* – Order of derivative in y direction.

4.2.49 splitTextLines Operation

Expects input images to have horizontal baselines in the center of the image. Following parameters are available:

- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- maxCut** *Float, Default: 2.0* – Maximal sum at split borders
- minWidth** *Float, Default: 4.0* – Minimal width of a child relative to image height.

4.2.50 subtractMean Operation

Mean normalization (Subtract mean) Following parameters are available:

- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file.

Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

4.2.51 tee Operation

tee writes the current set of images to the file system similarly to the linux command tee. The images in the pipeline remain unchanged. Following parameters are available:

- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- normalize** *Integer, Default: 1* – Normalize between 0 and 255 before storing.
- outputPath** *String, Default: tee-%base%idx%ext* – Path to the output files. See the -outputPath parameter of preprocessor for more information.

4.2.52 textSkewCorrection Operation

Corrects the text skew resulting from italic writing styles. Assumes that the base line is centered, i.e. pixels on the middle horizontal lines are not affected by this transformation. Following parameters are available:

- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- fromDegree** *Float, Default: -50.0* – Minimal degree considered by Hough transform.
- maxDegree** *Float, Default: 15.0* – If the detected text skew exceeds this parameter (in degrees), the detection is assumed to be incorrect and no correction is applied.
- resolution** *Integer, Default: 90* – Resolution of Hough transform.
- sobelKSize** *Integer, Default: 3* – Size of the sobel kernel.
- toDegree** *Float, Default: 50.0* – Maximal degree considered by Hough transform.

4.2.53 thinning Operation

Line thinning as proposed by T.Y. Zhang and C.Y. Suen. Following parameters are available:

- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

4.2.54 threshold Operation

The threshold command creates binary images. It is based on the OpenCV function `threshold()` and thus supports simple, adaptive, and otsu's thresholding. Following parameters are available:

<code>configDumpFileName</code>	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
<code>configFile</code>	<i>String, Default: <not set></i> – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
<code>maxVal</code>	<i>Integer, Default: 255</i> – Positive value after binarization (passed through to OpenCV).
<code>threshold</code>	<i>Float, Default: 127.0</i> – Binarization threshold (passed through to OpenCV).
<code>type</code>	<i>String, Default: BINARY,OTSU</i> – Thresholding type. See OpenCV's documentation for the threshold function. Connect options with ','. Available options are: BINARY, BINARY_INV, TRUNC, TOZERO, TOZERO_INV, OTSU

4.2.55 transpose Operation

Transpose the image. Following parameters are available:

<code>configDumpFileName</code>	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
<code>configFile</code>	<i>String, Default: <not set></i> – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.

4.2.56 vertTextSegmentation Operation

Extracts vertical cuts of text areas. The areas are identified by a constant horizontal gradient of the vertical projection. Following parameters are available:

<code>concatChildren</code>	<i>Integer, Default: 0</i> – Set to 0 to pass each found text segment separately through the pipeline. Set to 1 to concat all children to a single image removing the non-text areas.
<code>configDumpFileName</code>	<i>String, Default: <not set></i> – This can be used to write a file containing all parameters of the used configuration.
<code>configFile</code>	<i>String, Default: <not set></i> – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
<code>dilate</code>	<i>Float, Default: 0.05</i> – Safety margin which is added to the found text segments, relative to the page width. Set to negative to disable dilation.
<code>minMargin</code>	<i>Float, Default: 0.05</i> – Minimal vertical distance between two text areas relative to the page width.
<code>minSlope</code>	<i>Float, Default: 0.25</i> – Minimal gradient in the vertical projection of a

text area. The values of the vertical projection range from 0 to 1, and the projection is resized to -resolution

- minWidth** *Float, Default: 0.1* – Minimal width of a text area relative to the page width.
- morph** *String, Default: openFirst* – Morphology operations on the segmentation.
- openFirst: opening, then closing
 - closeFirst: closing, then opening
 - none: No morphology operation.
- outputSegmentation** *Integer, Default: 0* – Output an 1xresolution array indicating the classification of the columns in text and non-text columns. If this is set to 0, the input image is cutted according the text segmentations.
- resolution** *Integer, Default: 100* – Resolution for the vertical projection. 100 is a good value even for documents with largely differing sizes.
- transpose** *Integer, Default: 0* – Set to 1 to transpose the image before applying segmentation algorithm. This results producing horizontal instead of vertical cuts.

4.2.57 writeRects Operation

Writes the rectangles for each child into a file. Following parameters are available:

- configDumpFileName** *String, Default: <not set>* – This can be used to write a file containing all parameters of the used configuration.
- configFile** *String, Default: <not set>* – Configuration file (format: <key> <val>). Parameters in this file override values in the default configuration file. Command line arguments override all other settings. If this parameter is used within a config file, it has 'include once' semantics.
- outputPath** *String, Default: %base%idx-rect.txt* – Path to the rectangle files to create. The same placeholders as in the global outputPath can be used.

References

- G. Bradski and A. Kaehler. *Learning OpenCV: Computer vision with the OpenCV library*. O'Reilly Media, Inc., 2008.
- X. Huang, A. Acero, H.-W. Hon, and Raj R. *Spoken language processing: A guide to theory, algorithm, and system development*. Prentice Hall PTR, 2001.
- D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, N. Goel, M. Hannemann, P. Motlíček, Y. Qian, P. Schwarz, et al. The Kaldi speech recognition toolkit. 2011.
- J. A. Sanchez, A. H. Toselli, V. Romero, and E. Vidal. ICDAR2015 Competition HTRtS: Handwritten Text Recognition on the tranScriptorium Dataset. In *ICDAR*. IEEE, 2015.
- F. Stahlberg and S. Vogel. Document Skew Detection Based on Hough Space Derivatives. In *ICDAR*. IEEE, 2015a.
- F. Stahlberg and S. Vogel. Detecting dense foreground stripes in Arabic handwriting for accurate baseline positioning. In *ICDAR*. IEEE, 2015b.

Index

- baseline estimation, 14
- batch mode, 5, 8

- code, 6, 7, 10

- document skew detection, 13

- efficiency, 5

- feature extraction, 13

- global parameters, 15
 - configDumpFileName, 15
 - configFile, 16
 - idLength, 16
 - inputFile, 7, 16
 - logLevel, 16
 - nThreads, 16
 - outputPath, 7, 16
 - pipeline, 17
 - silentOverwrite, 18
 - singlePopulation, 18

- ICDAR2015 Competition HTRtS, 11

- Java, 5
- JavaDoc, 7

- Kaldi, 4, 9

- memory, 15
- modularity, 5

- NoClassDefFoundError, 15

- OpenCV, 5, 5, 6, 15
- OPENCV_JAR_PATH, 6
- OPENCV_NATIVE_LIB, 6, 15
- operation, 10, 18
 - adaptiveThreshold, 18
 - axisAlignedHough, 19
 - bbq, 19
 - blur, 20
 - col2graph, 20
 - componentDensity, 20
 - concat, 20
 - convertToFloat, 21
 - cutWithAltecXml, 21
 - cutWithPageXml, 21
 - devNull, 22
 - drawChildren, 22
 - drawKaldiAlignment, 22
 - drawTextLines, 23
 - exactOrientationCorrection, 23
 - extend, 24
 - extendForHoughsquare, 24
 - extractConstantRegions, 24
 - featExtract, 9, 24
 - fillTransparency, 26
 - filter, 26
 - flip, 27
 - grayscale, 8, 27
 - hough, 27
 - houghTextLine, 28
 - invert, 29
 - log, 29
 - morph, 29
 - multiChannelOtsu, 30
 - normalize, 30
 - normalizeText, 31
 - normalizeUpperBaseline, 31
 - orientationCorrection, 32
 - outlierRemove, 32
 - polynomialTextLine, 33
 - printMax, 33
 - projectionLineSegmentation, 33
 - rectSum, 34
 - reduce, 34
 - reducedAlcmTransform, 35
 - removeDiacritics, 36
 - removeLargeComponents, 36
 - removeSmallComponents, 36
 - removeUnderline, 37
 - removeVertTextMargin, 37
 - renderPageXmlTranscriptions, 35
 - scale, 8, 38
 - sobel, 38
 - splitTextLines, 38
 - subtractMean, 38
 - tee, 39
 - textSkewCorrection, 39
 - thinning, 39
 - threshold, 40
 - transpose, 8, 40
 - vertTextSegmentation, 40
 - writeRects, 41
- operation parameters, 18
 - aboveBaseline, 31
 - adaptiveType, 18
 - align, 35
 - alignmentFile, 22
 - analysisMode, 34
 - angleResolution, 27
 - angleSamplingFactor, 27
 - baselineHeight, 25
 - belowBaseline, 31
 - bgColorEstimateBorder, 35
 - blackDiscount, 30
 - blockSize, 18
 - blurMode, 28

blurRho, 28
 blurTheta, 28
 border, 21
 borderFactor, 28
 borderHeight, 22
 borderType, 38
 bottom, 24
 C, 18
 center, 20
 closeSize, 24
 combination, 28
 concatChildren, 40
 configDumpFileName, 19–41
 configFile, 19–41
 criterion, 23
 cutLevel, 21
 dataPoints, 33
 delayDelta, 25
 delayRaw, 25
 delays, 25
 deleteAboveAscenders, 28
 deleteBelowDescenders, 28
 deltaLambda, 28
 dilate, 40
 dim, 35
 drawSkewLine, 23
 endLambdaBandMin, 28
 endLambdaBandWidth, 28
 eps, 23
 extractors, 9, 25
 extractRegions, 21, 30
 extractTextObjects, 21, 30
 featAnhdfConnecticityTolerance, 25
 featAnhdfReductionMode, 25
 featAnhdfSegmentNum, 25
 featConcavityBaselineDependence, 25
 featDirectionalRadius, 25
 featRawCellHeight, 25
 featRawCellShift, 26
 featRawCellWidth, 26
 featRunlengthsNonNegative, 26
 featRunlengthsRadius, 26
 featSnakeAddCenterDistances, 26
 featSnakeAddRelativeFeats, 26
 featSnakeBackground, 26
 featSnakeDefaultHeight, 26
 featSnakeForeground, 26
 featSnakeNumber, 26
 fixX, 34
 fixY, 34
 flipCode, 27
 fontFamily, 35
 foregroundThreshold, 26, 37
 fromDegree, 39
 fromTheta, 19
 graphWidth, 20
 horizWeight, 23
 houghLineMode, 23
 houghMax, 28
 kaldiFile, 26
 kaldiId, 22, 26
 keepCoreZoneAspectRatio, 31
 kernelShape, 30
 kernelSize, 30
 kHeight, 35
 kSize, 38
 kWidth, 35
 left, 24
 maxAngle, 23, 32
 maxBelowShrink, 31
 maxBelowStretch, 31
 maxCut, 31, 32, 38
 maxDegree, 28, 39
 maxFontSize, 35
 maxForegroundFraction, 30
 maxHeight, 34, 36, 37
 maxIter, 30
 maxLineHeight, 34
 maxLineHeightVariance, 34
 maxProjectionRatio, 34
 maxRelHeight, 37
 maxStretchFactor, 32
 maxThickness, 37
 maxVal, 19, 40
 maxValArea, 28
 maxVariation, 37
 maxWidth, 34, 36
 minAspectRatio, 26
 minCroppedAboveAbsolute, 31, 32
 minCroppedAboveRatio, 31, 32
 minCroppedBelowRatio, 31
 minFontSize, 36
 minHeight, 27, 33, 37
 minLength, 24
 minLevel, 22, 30, 36
 minLineCount, 34
 minLineHeight, 34
 minMargin, 40
 minRelWidth, 37
 minRowSum, 37
 minSize, 20
 minSlope, 40
 minStretchFactor, 32
 minTextHeight, 37
 minWidth, 37, 38, 41
 mode, 20, 35
 morph, 41
 newMax, 30
 newMin, 31
 newUpperBaseline, 32
 noCorrection, 23, 32
 normalize, 39

- normalizeRegionChannels, 30
- noTextLineOperation, 28
- offset, 22
- openSize, 24
- operation, 29, 30, 32, 33
- order, 33
- outlierFactor, 33
- outputPath, 39, 41
- outputSegmentation, 41
- refine, 23
- reloadOriginal, 23, 32
- renderWithOpenCV, 36
- resolution, 23, 29, 32, 35, 39, 41
- right, 24
- sinusExp, 34
- sobelKSize, 24, 32, 39
- startLambdaBandMin, 29
- startLambdaBandWidth, 29
- thetaResolution, 19
- thickness, 22
- thicknessFactor, 37
- threshold, 19, 24, 33, 40
- toDegree, 39
- tolerance, 33
- toMatrix, 27
- top, 24
- toTheta, 19
- transpose, 22, 41
- truPath, 29
- type, 19, 40
- upperBaselineHighest, 32
- upperBaselineLowest, 32
- useIndexAttributes, 21
- usePageIds, 22
- useTruIfAvailable, 29
- winShift, 26
- winWidth, 26
- xmlPath, 21, 22, 30, 36
- xOrder, 38
- xScale, 38
- xSize, 20
- yOrder, 38
- yScale, 38
- ySize, 20

parameters

- global, *see* global parameters
- operation, *see* operation parameters

pipeline, 5, 11

repository, 10

scalability, 5

UnsatisfiedLinkError, 15