

PooledTimeSeriesParser

The [Pooled Time Series algorithm](#) was developed by [Michael Ryoo](#) and it allows for video descriptors to be considered over time and in this consideration for videos to be compared based on the activity going on in their scenes. In short, Pooled Time Series is a video comparison metric. An ALv2 licensed version of the [Pooled Time Series code](#) is available for use in computing Histogram of Oriented Gradients (HOG) and Histogram of Optical Flows (HOF) which can be useful extracted data and metadata for a Tika Parser.

Metadata Representation

The ultimate goal of the project is to be able to extract metadata and data from videos and to index that information inside of a search engine like Apache Solr. Videos, like images, are just numbers - or a ordered sequence of number - or matrices. There are many ways in which these numbers can be defined. Some popular visual descriptors are Histogram of Gradients, Optical Flow vectors, RGB or Color Histograms. The challenge is to figure out a way to map this datatype to a datatype that can be understood by Solr. In the case of color based histograms, we can convert the image into a matrix of hex values, where each hex value is the pixel color value and index that as a text_ws field in Solr.

Some Related Efforts

[Shutterstock](#) developed an [image search tool](#) using a similar approach.

Representation of output data

The data output from the Pooled Time Series parser is an XHTML document of table values, where each `<tr>..</tr>` would be a row of the feature matrix and `<td>` would be the corresponding element in that column. When using a search engine like Apache Solr to do ranking or querying we can to compute a distance function on these values (for the dataset and the query video), such as Chi-Squared, which is what the pooled time series algorithm does.

A Tika Parser has been developed that implements the Pooled Time Series algorithm above and that outputs the HOF and HOG data from videos for use in later processing and indexing. Read on below to install and use it!

Pre-requisites

Install Pooled Time Series

1. `mkdir -p $HOME/git && cd $HOME/git && git clone https://github.com/USCDataScience/hadoop-pot.git`
2. `cd hadoop-pot/hadoop-pot-assembly && mvn install assembly:assembly`
3. Follow steps 3, 4 and 5 from [the install guide from Pooled Time Series](#) and confirm that `pooled_time_series` installed correctly. Note the [pre-requisites from Pooled Time Series](#) require you to install OpenCV and set some environment variables.

After above steps you must be able to execute `pooled_time_series` through terminal and get below output

```
usage: pooled_time_series
-d,--dir <directory>      A directory with image files in it
-f,--file <file>          Path to a single file
-h,--help                  Print this message.
-j,--json                  Set similarity output format to JSON.
                           Defaults to .txt
-o,--outputfile <output file> File containing similarity results.
                           Defaults to ./similarity.txt
-p,--pathfile <path file>  A file containing full absolute paths to
                           videos. Previous default was
                           memex-index_temp.txt
```

Run PooledTimeSeries Parser from Tika-App

Grab an MP4, [QuickTime](#), or other video (supported by your OpenCV implementation). Then run the following command:

```
java -classpath tika-app/target/tika-app-X.Y.jar org.apache.tika.cli.TikaCLI --config=$HOME/git/pooled_time_series/src/main/resources/tika-config.xml -J yourfile.mov
```

which should output something like:

[illegible]

Will this work from Tika Server?

Yes, it will! Start Tika server like so:

```
java -jar tika-server/target/tika-server-X.Y.jar --config=$HOME/git/pooled_time_series/src/main/resources/tika-config.xml
```

Then send your movie file to Tika server like so:

```
curl -T yourfile.mov http://localhost:9998/rmeta
```

which should produce in response:

[illegible]