

# TikaAndNLTK

## Natural Language Toolkit (NLTK) with Tika

NLTK is a python library and in order to integrate and use its capabilities with Tika one must install the server endpoint created for NLTK using Flask to extract Named Entities from text passed to it.

- [Natural Language Toolkit \(NLTK\) with Tika](#)
  - [Installation](#)
  - [Start NLTK Server](#)
  - [Preparing resources for NLTK in Tika-App](#)
  - [Running NLTK with Tika](#)

### Installation

#### 1. Simple with pip

```
$ pip install --process-dependency-links nltkrest
```

2. **Setuptools and/or Distribute:** The module can be downloaded from [github](#) and then installed with the following commands:

```
$ cd NLTKRest/nltkrest
$ python setup.py install nltkrest
```

After installation you will have a command called *nltk-server*. By default it starts a REST server on port 8881. You can change the port by typing *--port* or *-p*. You can also turn on verbose mode by typing *-v*.

### Start NLTK Server

```
$ nltk-server -v --port 8888
```

You should see a message:

Running on <http://127.0.0.1:8888> The server is up and ready!

### Preparing resources for NLTK in Tika-App

You can either perform steps 1 & 2 together or just 3.

1. **Activate Named Entity Parser** In order to use any of the [NamedEntityParser](#) implementations in Tika , the parser responsible for handling the name recognition task needs to be enabled. This can be done with Tika Config XML file, as follows

```
<?xml version="1.0" encoding="UTF-8"?>
<properties>
  <parsers>
    <parser class="org.apache.tika.parser.ner.NamedEntityParser">
      <mime>text/plain</mime>
      <mime>text/html</mime>
      <mime>application/xhtml+xml</mime>
    </parser>
  </parsers>
</properties>
```

This configuration has to be supplied in the later phases, so store it as 'tika-config.xml'.

#### 2. Supply NLTKServer.properties file

It is imperative that Tika should know on what host you are running the *nltk-server*. By default Tika will assume your server runs on port 8881. In order to specify any other port, you must supply a NLTKServer.properties file. Sample NLTKServer.properties file. My file looks like the following:

```
nltk.server.url=http://localhost:8881
```

In an nutshell

```
#Create a directory for keeping the config and properties file.
export NLTK_RES=$HOME/NLTKRest-resources
mkdir -p $NLTK_RES
cd $NLTK_RES
#config file must be stored in this directory
pwd

export PATH_PREFIX="$NLTK_RES/org/apache/tika/parser/ner/nltk"
mkdir -p $PATH_PREFIX
#create and edit the properties file
vim $PATH_PREFIX/NLTKServer.properties
```

### 3. Download NLTKRest-resources

Better yet, you could skip the previous two steps completely and save the hassle. Simply download the project [NLTKRest-resources](#) and edit the properties file

```
cd $HOME && git clone https://github.com/manalishah/NLTKRest-resources
export NLTK_RES=$HOME/NLTKRest-resources
vim $NLTK_RES/org/apache/tika/parser/ner/nltk/NLTKServer.properties
```

## Running NLTK with Tika

Finally, we've reached this point where we can smile and let Tika do the working!

```
export TIKA_APP={your/path/to/tika-app}/target/tika-app-1.13-SNAPSHOT.jar

#set the system property to use NLTKNERecogniser class
java -Dner.impl.class=org.apache.tika.parser.ner.nltk.NLTKNERecogniser -classpath $NLTK_RES:$TIKA_APP org.apache.tika.cli.TikaCLI --config=$NLTK_RES/tika-config.xml -m http://www.hawking.org.uk/
```

This will output metadata along with named entities extracted using nltk:

NER\_NAMES: Gonville  
NER\_NAMES: Einstein  
NER\_NAMES: Briefer History  
NER\_NAMES: ALS  
NER\_NAMES: Unbreakable Code  
NER\_NAMES: Click  
NER\_NAMES: Cambridge  
NER\_NAMES: George  
NER\_NAMES: Lucasian Professor  
NER\_NAMES: Stephen Hawking  
NER\_NAMES: Latest  
NER\_NAMES: Hubble  
NER\_NAMES: 1979  
NER\_NAMES: Brief History  
NER\_NAMES: California  
NER\_NAMES: University  
NER\_NAMES: 1663  
NER\_NAMES: 1982  
NER\_NAMES: London  
NER\_NAMES: US National Academy  
NER\_NAMES: Baby Universe  
NER\_NAMES: Home About Stephen The Computer Stephen  
NER\_NAMES: Applied Mathematics  
NER\_NAMES: Santa Cruz  
NER\_NAMES: Leiden University  
NER\_NAMES: CBE  
NER\_NAMES: Science  
NER\_NAMES: Caius College  
NER\_NAMES: HUDF09 Team  
NER\_NAMES: Dennis Stanton Avery  
NER\_NAMES: 2009  
NER\_NAMES: ESA  
NER\_NAMES: Annie  
NER\_NAMES: NASA  
NER\_NAMES: Black Holes  
NER\_NAMES: Universe  
NER\_NAMES: Sally Tsui  
NER\_NAMES: eXtreme Deep Field  
NER\_NAMES: Centre  
NER\_NAMES: Royal Society  
NER\_NAMES: Weebly  
NER\_NAMES: Cambridge Lectures Publications Books Images Films Videos Stephen  
NER\_NAMES: 1963  
NER\_NAMES: Unbreakable Code Hot  
NER\_NAMES: Isaac Newton  
NER\_NAMES: Lucy  
NER\_NAMES: Stephen