

---

# **clstk Documentation**

**Nisarg Jhaveri**

**Jun 02, 2018**



---

## Getting Started

---

<b>1 Installation</b>	<b>3</b>
1.1 Clone repository . . . . .	3
1.2 Python dependencies . . . . .	3
1.3 Setup CLUTO (optional) . . . . .	3
1.4 Setup ROUGE 1.5.5 (optional) . . . . .	3
1.5 Setup dependencies for TQE (optional) . . . . .	4
1.6 Setup NeuralTextSimplification (optional) . . . . .	4
<b>2 Datasets</b>	<b>5</b>
2.1 DUC 2004 Gujarati Dataset . . . . .	5
2.2 MultiLing Pilot 2011 Dataset . . . . .	5
<b>3 Usage</b>	<b>7</b>
3.1 Summarize . . . . .	7
3.2 Evaluate . . . . .	8
<b>4 Core</b>	<b>9</b>
4.1 Sentence class . . . . .	9
4.2 SentenceCollection class . . . . .	10
4.3 Corpus class . . . . .	12
4.4 Summary class . . . . .	12
<b>5 Utils</b>	<b>13</b>
5.1 nlp utils . . . . .	13
5.2 ProgressBar class . . . . .	14
<b>6 Evaluate</b>	<b>15</b>
6.1 RougeScore . . . . .	15
6.2 ExternalRougeScore . . . . .	16
<b>7 Translate</b>	<b>17</b>
7.1 googleTranslate . . . . .	17
7.2 googleTranslateWeb . . . . .	18
<b>8 Simplify</b>	<b>19</b>
8.1 neuralTextSimplification . . . . .	19
<b>9 Translation Quality Estimation</b>	<b>21</b>

---

9.1 qualityEstimation . . . . .	21
<b>Python Module Index</b>	<b>23</b>

`clstk` is a free and open source tool-kit for cross-lingual summarization (CLS). The tool-kit is intended for the use of both, developers and researchers working on cross-lingual summarization. End-users wanting to use cross-lingual summarization in real-world applications can also benefit from the tool-kit.

The tool-kit currently contains implementation of three CLS methods along with bootstrap code and other modules required for CLS. We encourage developers to contribute more methods in the tool-kit.



# CHAPTER 1

---

## Installation

---

### 1.1 Clone repository

```
$ git clone https://github.com/nisargjhaveri/clstk
```

### 1.2 Python dependencies

The dependencies are listed in `requirements.txt`.

To install all the dependencies, run `pip` as followed.

```
$ pip install --upgrade -r requirements.txt
```

Also install `nltk` packages called `stopwords` and `punkt`.

```
$ python -m nltk.download stopwords punkt -d $NLTK_DATA
```

### 1.3 Setup CLUTO (optional)

<http://glaros.dtc.umn.edu/gkhome/cluto/cluto/download>

This is required if you want to use “linBilmes” summarizer.

Set an environment variable `CLUTO_BIN_PATH` with the path of directory containing `vcluster` binary file.

### 1.4 Setup ROUGE 1.5.5 (optional)

<https://github.com/nisargjhaveri/ROUGE-1.5.5-unicode>

This is required only if you plan to evaluate the summaries using ROUGE score.

Obtain and setup ROUGE 1.5.5 according to the instructions there.

Set an environment variable ROUGE\_HOME with the path to ROUGE root directory, the one containing ROUGE-1.5.5.pl file.

## **1.5 Setup dependencies for TQE (optional)**

<https://github.com/nisargjhaveri/tqe>

Install dependencies for tqe module according to the details provided in the link above.

## **1.6 Setup NeuralTextSimplification (optional)**

<https://github.com/senisioi/NeuralTextSimplification>

Setup system from above URL and set NTS\_OPENNMT\_PATH, NTS\_MODEL\_PATH and NTS\_GPUS variables accordingly.

# CHAPTER 2

---

## Datasets

---

### 2.1 DUC 2004 Gujarati Dataset

<https://github.com/nisargjhaveri/duc2004-translated>

This is a cross-lingual summarization evaluation dataset for English to Gujarati summarization. The dataset can be obtained from the link mentioned above.

You'll also need original DUC 2004 dataset as the above link does not contain source documents due to licensing reasons.

### 2.2 MultiLing Pilot 2011 Dataset

<http://users.iit.demokritos.gr/~ggianna/TAC2011/MultiLing2011.html>

This dataset contains parallel document sets in seven languages: English, Arabic, Czech, French, Greek, Hebrew and Hindi. Summaries for each document set is available in all languages, making the dataset usable for cross-lingual summarization evaluation.

The data needs to be cleaned and formatted for the use with `clstk`.



# CHAPTER 3

---

## Usage

---

### 3.1 Summarize

`sum.py` is used to summarize a document set. It allows selecting specific CLS method and parameters for that particular method.

```
$ python sum.py --help
usage: sum.py [-h] {linBilmes,coRank,simFusion} ...
Automatically summarize a set of documents

optional arguments:
  -h, --help            show this help message and exit

methods:
  Summarization method

  {linBilmes,coRank,simFusion}
```

The following command shows help for selected method.

```
$ python sum.py {method} --help
```

Following is the common pattern to run a CLS method on one document set.

```
$ python sum.py {method} [options] {source_directory}
```

All files stored in the directory `source_directory` are read and treated as a part of document set to summarize. The files are expected to be plain text files.

#### 3.1.1 Required arguments

**source\_directory** Directory containing a set of files to be summarized.

### 3.1.2 Common options

Here is a list of common optional arguments across all CLS methods.

<b>-h, --help</b>	show this help message and exit
<b>-v, --verbose</b>	Show verbose information messages
<b>--no-colors</b>	Don't show colors in verbose log
<b>-s N, --size N</b>	Maximum size of the summary
<b>-w, --words</b>	Caluated size as number of words instead of characters
<b>--source-lang lang</b>	Two-letter language code of the source documents language. Defaults to <i>en</i>
<b>-l lang, --target-lang lang</b>	Two-letter language code to generate cross-lingual summary. Defaults to source language.

## 3.2 Evaluate

Another script called `evaluate.py` is used to run and evaluate CLS methods over a CLS evaluation dataset.

Similar to `sum.py`, this script also needs the CLS method as first argument and other argument follows depending on the selected method.

```
$ python evaluate.py {method} [options] {source_path} {models_path} {summaries_path}
```

### 3.2.1 Required arguments

**source\_path** Directory containing all the source files to be summarized. Each set of documents are expected to be in different directories inside this path.

**models\_path** Directory containing all the model summaries. Each set of summaries are expected to be in different directory inside this path, having the same name as the corresponding directory in the source directory.

**summaries\_path** Directory to store the generated summaries. The directory will be created if not already exists.

### 3.2.2 Common options

<b>-h, --help</b>	show this help message and exit
<b>--only-rouge</b>	Do not run summarizer. Only compile ROUGE score for existing summaries in summaries_path
<b>-s N, --size N</b>	Maximum size of the summary
<b>-w, --words</b>	Caluated size as number of words instead of characters
<b>--source-lang lang</b>	Two-letter language code of the source documents language. Defaults to <i>en</i>
<b>-l lang, --target-lang lang</b>	Two-letter language code to generate cross-lingual summary. Defaults to source language.

# CHAPTER 4

---

## Core

---

The core contains the bootstrap code for summarization needs. The core provides:

- A common standard structure for documents and summaries to ensure interoperability between different components.
- Utilities for loading document sets into the common structure.
- Common utilities on document sets, documents and sentences, for example sentence splitting, tokenization, etc.

### 4.1 Sentence class

```
class clstk.sentence.Sentence (sentenceText)
Bases: object

Class to represent a single sentence

__init__ (sentenceText)
    Set sentence text and translated text

    Parameters sentenceText – sentence text

setText (sentenceText)
    Set text for the sentence

    Parameters sentenceText – sentence text

getText ()
    Get sentence text

    Returns sentence text

setTranslation (translation)
    Set translated text

    Parameters translation – translated text
```

```
getTranslation()
    Get translated text

    The translated text defaults to sentence text

    Returns translated text

setVector (vector)
    Set sentence vector

    Parameters vector – sentence vector

getVector ()
    Get sentence vector

    Returns sentence vector

setTranslationVector (vector)
    Set sentence vector for translated text

    Parameters vector – sentence vector

getTranslationVector ()
    Get sentence vector for translated text

    Returns sentence vector

setExtra (key, value)
    Set extra key-value pair

    Parameters

        • key – key for the stored value
        • value – value to store

getExtra (key, default=None)
    Get extra value from key

    Parameters

        • key – key for the stored value
        • default – default value if key not found

charCount ()
    Get character count for translated text

    Returns Number of character in translated text

tokenCount ()
    Get token count for translated text

    Returns Number of tokens in translated text

__weakref__
    list of weak references to the object (if defined)
```

## 4.2 SentenceCollection class

```
class cltk.sentenceCollection.SentenceCollection
Bases: object

Class to store a colelction of sentences.
```

Also provides several common operations on the collection.

### `__init__()`

Initialize the collection

### `setSourceLang(lang)`

Set source language for the collection

**Parameters** `lang` – two-letter code for source language

### `setTargetLang(lang)`

Set target language for the collection

**Parameters** `lang` – two-letter code for target language

### `addSentence(sentence)`

Add a sentence to the collection

**Parameters** `sentence` – sentence to be added

### `addSentences(sentences)`

Add sentences to the collection

**Parameters** `sentences` – list of sentence to be added

**See also:**

`clstk.sentenceCollection.SentenceCollection.addSentence()`

### `getSentences()`

Get list of sentences in the collection

**Returns** list of sentences

### `getSentenceVectors()`

Get list of sentence vectors for sentences in the collection

**Returns** np.array containing sentence vectors

### `getTranslationSentenceVectors()`

Get list of sentence vectors for translations of sentences in the collection

**Returns** np.array containing sentence vectors

### `generateSentenceVectors()`

Generate sentence vectors

### `generateTranslationSentenceVectors()`

Generate sentence vectors for translations

### `translate(sourceLang, targetLang, replaceOriginal=False)`

Translate sentences

#### **Parameters**

- `sourceLang` – two-letter code for source language
- `targetLang` – two-letter code for target language
- `replaceOriginal` – Replace source text with translation if True. Used for early-translation

### `simplify(sourceLang, replaceOriginal=False)`

Simplify sentences

#### **Parameters**

- **sourceLang** – two-letter code for language
- **replaceOriginal** – Replace source sentences with simplified sentences. Used for early-simplify.

weakref

list of weak references to the object (if defined)

## 4.3 Corpus class

**class** clstk.corpus.Corporus (*dirname*)

Bases: clstk.sentenceCollection.SentenceCollection

Class for source documents. Contains utilities for loading document set.

\_\_init\_\_ (*dirname*)

Initialize the class

**Parameters** **dirname** – Directory from where source documents are to be loaded

**load** (*params*, *translate=False*, *replaceWithTranslation=False*, *simplify=False*, *replaceWithSimplified=False*)

Load source document set

**Parameters**

- **params** – dict containing different params including sourceLang and targetLang.
- **translate** – Whether to translate sentences to target language
- **replaceWithTranslation** – Whether to replace source sentences with translation
- **simplify** – Whether to simplify sentences
- **replaceWithSimplified** – Whether to replace source sentences with simplified sentences

## 4.4 Summary class

**class** clstk.summary.Summary

Bases: clstk.sentenceCollection.SentenceCollection

**charCount()**

Get total number of character in all the sentences

**tokenCount()**

Get total number of tokens in all the sentences

**getSummary()**

Get printable summary generated from source text

**getTargetSummary()**

Get printable summary generated from translated text

# CHAPTER 5

---

## Utils

---

Different utilities

### 5.1 nlp utils

`clstk.utils.nlp.getSentenceSplitter()`

Get sentence splitter function

**Returns** A function which takes a string and return list of sentence as strings.

`clstk.utils.nlp.getTokenizer(lang)`

Get tokenizer for a given language

**Parameters** `lang` – language

**Returns** tokenizer, which takes a sentence as string and returns list of tokens

`clstk.utils.nlp.getDetokenizer(lang)`

Get detokenizer for a given language

**Parameters** `lang` – language

**Returns** detokenizer, which takes list of tokens and returns a sentence as string

`clstk.utils.nlp.getStemmer()`

Get stemmer. For now returns Porter Stemmer

**Returns** stemmer, which takes a token and returns its stem

`clstk.utils.nlp.getStopwords(lang)`

Get list of stopwords for a given language

**Parameters** `lang` – language

**Returns** list of stopwords including common puncuations

## 5.2 ProgressBar class

```
class clstk.utils.progress.ProgressBar(totalCount)
Bases: object
```

Class to manage and show pretty progress-bar in the console

```
__init__(totalCount)
    Initialize the progressbar
```

**Parameters** **totalCount** – Total items to be processed

```
done(doneCount)
    Move progressbar ahead
```

**Parameters** **doneCount** – Out of **totalCount**, this many have been processed

```
complete()
    Complete progress
```

```
__weakref__
    list of weak references to the object (if defined)
```

# CHAPTER 6

---

## Evaluate

---

Module containing ROUGE implementations.

### 6.1 RougeScore

Python implementation of ROUGE score.

Taken and adopted from:

- <https://github.com/miso-belica/sumy/blob/master/sumy/evaluation/rouge.py>
- <https://github.com/google/seq2seq/blob/master/seq2seq/metrics/rouge.py>

```
class clstk.evaluation.rougeScore(tokenizer=None, stemmer=None)
Bases: object
```

Implementation of ROUGE score.

```
__init__(tokenizer=None, stemmer=None)
x.__init__(...) initializes x; see help(type(x)) for signature
```

```
rouge_n(summary, model_summaries, n=2)
Computes ROUGE-N of two text collections of sentences.
```

```
rouge_l_sentence_level(evaluated_sentences, reference_sentences)
Computes ROUGE-L (sentence level) of two text collections of sentences.
```

```
rouge_l_summary_level(evaluated_sentences, reference_sentences)
Computes ROUGE-L (summary level) of two text collections of sentences.
```

```
rouge(hyp_refs_pairs, print_all=False)
Calculates and prints average rouge scores for a list of hypotheses and references
```

#### Parameters

- **hyp\_refs\_pairs** – List containing pairs of path to summary and list of paths to reference summaries

- **print\_all** – Print every evaluation along with averages

**\_\_weakref\_\_**

list of weak references to the object (if defined)

## 6.2 ExternalRougeScore

Integration with external ROUGE tool-kit.

We recommend the use of <https://github.com/nisargjhaveri/ROUGE-1.5.5-unicode>  
ROUGE\_HOME variable needs to be set to run this.

**class** clstk.evaluation.externalRougeScore.**ExternalRougeScore**  
Bases: object

Integration with external ROUGE tool-kit.

**rouge** (*summaryRefsList*)

Runs external ROUGE-1.5.5 and prints results

**Parameters** **summaryRefsList** – List containing pairs of path to summary and list of paths  
to reference summaries

**\_\_weakref\_\_**

list of weak references to the object (if defined)

# CHAPTER 7

---

## Translate

---

Implementation of different translators.

In general you should not need to use these directly.

**See also:**

`clstk.sentenceCollection.SentenceCollection.translate()` and `clstk.corpus.Corporus.load()`

### 7.1 googleTranslate

Translate using Google Translate.

To use this, environmental variable `GOOGLE_APPLICATION_CREDENTIALS` needs to be set with file containing your key for Google Cloud account.

See <https://cloud.google.com/translate/docs/reference/libraries>

`clstk.translate.googleTranslate.translate(text, sourceLang, targetLang)`

Translate text

#### Parameters

- `text` – Text, each line contains one sentence
- `sourceLang` – Two-letter code for source language
- `targetLang` – Two-letter code for target language

**Returns** translated text and list of translated sentences

**Return type** (translation, sentences)

## 7.2 googleTranslateWeb

**DO NOT** use this for commercial purposes

```
cstk.translate.googleTranslateWeb.translate(text, sourceLang, targetLang, sentencePer-  
Line=True)
```

Translate text

### Parameters

- **text** – Text, each line contains one sentence
- **sourceLang** – Two-letter code for source language
- **targetLang** – Two-letter code for target language

**Returns** translated text and list of translated sentences

**Return type** (translation, sentences)

# CHAPTER 8

---

## Simplify

---

Simplify sentences using different methods.

In general you should not need to use these directly.

**See also:**

`clstk.sentenceCollection.SentenceCollection.simplify()` and `clstk.corpus.Corporus.load()`

### 8.1 neuralTextSimplification

Neural Text Simplification.

You need to set NTS\_OPENNMT\_PATH, NTS\_MODEL\_PATH and NTS\_GPUS environmental variables to use this.

`clstk.simplify.neuralTextSimplification.simplify(sentences, lang)`  
Simplify sentences using NTS

#### Parameters

- **sentences** – List of sentence
- **lang** – Language of sentences

**Returns** List of simplified sentences



# CHAPTER 9

---

## Translation Quality Estimation

---

Estimate translation quality

### 9.1 qualityEstimation

Translation Quality Estimation

Setup dependencies for TQE to use this. <https://github.com/nisargjhaveri/tqe>

You also need to train model using the said tqe system.

`clstk.qualityEstimation.qualityEstimation.estimate(sentenceCollection, modelPath)`  
Estimate translation quality for each sentence in collection. It sets an extra value with key `qeScore` on each sentence.

**Parameters** `sentenceCollection` `(clstk.sentenceCollection.SentenceCollection)` – SentenceCollection to estimate quality

**See also:**

`clstk.sentence.Sentence.getExtra()`



---

## Python Module Index

---

### C

clstk.evaluation.externalRougeScore, [16](#)  
clstk.evaluation.rougeScore, [15](#)  
clstk.qualityEstimation.qualityEstimation,  
    [21](#)  
clstk.simplify.neuralTextSimplification,  
    [19](#)  
clstk.translate.googleTranslate, [17](#)  
clstk.translate.googleTranslateWeb, [18](#)  
clstk.utils.nlp, [13](#)



### Symbols

`__init__()` (clstk.corpus.Corpus method), 12  
`__init__()` (clstk.evaluation.rougeScore.RougeScore method), 15  
`__init__()` (clstk.sentence.Sentence method), 9  
`__init__()` (clstk.sentenceCollection.SentenceCollection method), 11  
`__init__()` (clstk.utils.progress.ProgressBar method), 14  
`__weakref__` (clstk.evaluation.externalRougeScore.ExternalRougeScore attribute), 16  
`__weakref__` (clstk.evaluation.rougeScore.RougeScore attribute), 16  
`__weakref__` (clstk.sentence.Sentence attribute), 10  
`__weakref__` (clstk.sentenceCollection.SentenceCollection attribute), 12  
`__weakref__` (clstk.utils.progress.ProgressBar attribute), 14

### A

`addSentence()` (clstk.sentenceCollection.SentenceCollection method), 11  
`addSentences()` (clstk.sentenceCollection.SentenceCollection method), 11

### C

`charCount()` (clstk.sentence.Sentence method), 10  
`charCount()` (clstk.summary.Summary method), 12  
clstk.evaluation.externalRougeScore (module), 16  
clstk.evaluation.rougeScore (module), 15  
clstk.qualityEstimation.qualityEstimation (module), 21  
clstk.simplify.neuralTextSimplification (module), 19  
clstk.translate.googleTranslate (module), 17  
clstk.translate.googleTranslateWeb (module), 18  
clstk.utils.nlp (module), 13  
`complete()` (clstk.utils.progress.ProgressBar method), 14  
Corpus (class in clstk.corpus), 12

### D

`done()` (clstk.utils.progress.ProgressBar method), 14

### E

`estimate()` (in module clstk.qualityEstimation.qualityEstimation), 21

ExternalRougeScore (class in clstk.evaluation.externalRougeScore), 16

### G

`generateSentenceVectors()` (clstk.sentenceCollection.SentenceCollection method), 11  
`generateTranslationSentenceVectors()` (clstk.sentenceCollection.SentenceCollection method), 11  
`getDetokenizer()` (in module clstk.utils.nlp), 13  
`getExtra()` (clstk.sentence.Sentence method), 10  
`getSentences()` (clstk.sentenceCollection.SentenceCollection method), 11  
`getSentenceSplitter()` (in module clstk.utils.nlp), 13  
`getSentenceVectors()` (clstk.sentenceCollection.SentenceCollection method), 11  
`getStemmer()` (in module clstk.utils.nlp), 13  
`getStopwords()` (in module clstk.utils.nlp), 13  
`getSummary()` (clstk.summary.Summary method), 12  
`getTargetSummary()` (clstk.summary.Summary method), 12

`getText()` (clstk.sentence.Sentence method), 9  
`getTokenizer()` (in module clstk.utils.nlp), 13  
`getTranslation()` (clstk.sentence.Sentence method), 9  
`getTranslationSentenceVectors()` (clstk.sentenceCollection.SentenceCollection method), 11

`getTranslationVector()` (clstk.sentence.Sentence method), 10

`getVector()` (clstk.sentence.Sentence method), 10

### L

`load()` (clstk.corpus.Corpus method), 12

## P

ProgressBar (class in clstk.utils.progress), [14](#)

## R

rouge() (clstk.evaluation.externalRougeScore.ExternalRougeScore  
method), [16](#)  
rouge() (clstk.evaluation.rougeScore.RougeScore  
method), [15](#)  
rouge\_1\_sentence\_level()  
(clstk.evaluation.rougeScore.RougeScore  
method), [15](#)  
rouge\_1\_summary\_level()  
(clstk.evaluation.rougeScore.RougeScore  
method), [15](#)  
rouge\_n() (clstk.evaluation.rougeScore.RougeScore  
method), [15](#)  
RougeScore (class in clstk.evaluation.rougeScore), [15](#)

## S

Sentence (class in clstk.sentence), [9](#)  
SentenceCollection (class in clstk.sentenceCollection), [10](#)  
setExtra() (clstk.sentence.Sentence method), [10](#)  
setSourceLang() (clstk.sentenceCollection.SentenceCollection  
method), [11](#)  
setTargetLang() (clstk.sentenceCollection.SentenceCollection  
method), [11](#)  
setText() (clstk.sentence.Sentence method), [9](#)  
setTranslation() (clstk.sentence.Sentence method), [9](#)  
setTranslationVector() (clstk.sentence.Sentence method),  
[10](#)  
setVector() (clstk.sentence.Sentence method), [10](#)  
simplify() (clstk.sentenceCollection.SentenceCollection  
method), [11](#)  
simplify() (in module clstk.simplify.neuralTextSimplification),  
[19](#)  
Summary (class in clstk.summary), [12](#)

## T

tokenCount() (clstk.sentence.Sentence method), [10](#)  
tokenCount() (clstk.summary.Summary method), [12](#)  
translate() (clstk.sentenceCollection.SentenceCollection  
method), [11](#)  
translate() (in module clstk.translate.googleTranslate), [17](#)  
translate() (in module  
clstk.translate.googleTranslateWeb), [18](#)