
gTTS Documentation

Pierre-Nick Durette

Jun 18, 2021

DOCUMENTATION

1	Installation	3
1.1	Command-line (<code>gtts-cli</code>)	3
1.1.1	<code>gtts-cli</code>	3
1.1.2	Examples	4
1.1.3	Playing sound directly	5
1.2	Module (<code>gtts</code>)	5
1.2.1	<code>gTTS</code> (<code>gtts.gTTS</code>)	5
1.2.2	Languages (<code>gtts.lang</code>)	7
1.2.3	Localized ‘accents’	7
1.2.4	Examples	8
1.2.5	Playing sound directly	8
1.2.6	Logging	8
1.3	Pre-processing and tokenizing	9
1.3.1	Definitions	9
1.3.2	Pre-processing	9
1.3.3	Tokenizing	11
1.3.4	Minimizing	12
1.3.5	<code>gtts.tokenizer</code> module reference (<code>gtts.tokenizer</code>)	12
1.4	License	16
1.5	Contributing	16
1.5.1	Reporting Issues	16
1.5.2	Submitting Patches	17
1.5.3	Testing	17
1.6	Changelog	17
1.6.1	2.2.3 (2021-06-17)	17
1.6.2	2.2.2 (2021-02-03)	18
1.6.3	2.2.1 (2020-11-15)	18
1.6.4	2.2.0 (2020-11-14)	18
1.6.5	2.1.2 (2020-11-10)	19
1.6.6	2.1.1 (2020-01-25)	19
1.6.7	2.1.0 (2020-01-01)	19
1.6.8	2.0.4 (2019-08-29)	20
1.6.9	2.0.3 (2018-12-15)	20
1.6.10	2.0.2 (2018-12-09)	21
1.6.11	2.0.1 (2018-06-20)	21
1.6.12	2.0.0 (2018-04-30)	21
1.6.13	1.2.2 (2017-08-15)	23
1.6.14	1.2.1 (2017-08-02)	24
1.6.15	1.2.0 (2017-04-15)	24
1.6.16	1.1.8 (2017-01-15)	24

1.6.17	1.1.7 (2016-12-14)	24
1.6.18	1.1.6 (2016-07-20)	25
1.6.19	1.1.5 (2016-05-13)	25
1.6.20	1.1.4 (2016-02-22)	25
1.6.21	1.1.3 (2016-01-24)	25
1.6.22	1.1.2 (2016-01-13)	25
1.6.23	1.0.7 (2015-10-07)	26
1.6.24	1.0.6 (2015-07-30)	26
1.6.25	1.0.5 (2015-07-15)	26
1.6.26	1.0.4 (2015-05-11)	26
1.6.27	1.0.3 (2014-11-21)	26
1.6.28	1.0.2 (2014-05-15)	27
1.6.29	1.0.1 (2014-05-15)	27
1.6.30	1.0 (2014-05-08)	27
2	Misc	29
	Python Module Index	31
	Index	33

gTTS (*Google Text-to-Speech*), a Python library and CLI tool to interface with Google Translate's text-to-speech API. Writes spoken mp3 data to a file, a file-like object (bytestring) for further audio manipulation, or `stdout`. It features flexible pre-processing and tokenizing.

INSTALLATION

```
pip install gTTS
```

1.1 Command-line (gtts-cli)

After installing the package, the `gtts-cli` tool becomes available:

```
$ gtts-cli
```

1.1.1 gtts-cli

Read `<text>` to mp3 format using Google Translate's Text-to-Speech API (set `<text>` or `-file <file>` to `-` for standard input)

```
gtts-cli [OPTIONS] <text>
```

Options

- f, --file <file>**
Read from `<file>` instead of `<text>`.
- o, --output <file>**
Write to `<file>` instead of stdout.
- s, --slow**
Read more slowly.
- l, --lang <lang>**
IETF language tag. Language to speak in. List documented tags with `-all`.
Default en
- t, --tld <tld>**
Top-level domain for the Google host, i.e <https://translate.google.<tld>>
Default com
- nocheck**
Disable strict IETF language tag checking. Allow undocumented tags.

--all

Print all documented available IETF language tags and exit.

--debug

Show debug information.

--version

Show the version and exit.

Arguments

<text>

Optional argument

1.1.2 Examples

List available languages:

```
$ gtts-cli --all
```

Read 'hello' to hello.mp3:

```
$ gtts-cli 'hello' --output hello.mp3
```

Read "c'est la vie" in French to cestlavie.mp3:

```
$ gtts-cli "c'est la vie" --lang fr --output cestlavie.mp3
```

Read '' to .mp3 (in Mandarin, using google.cn):

```
$ gtts-cli '' --tld cn --lang zh-cn --output .mp3
```

Read 'slow' slowly to slow.mp3:

```
$ gtts-cli 'slow' --slow --output slow.mp3
```

Read 'hello' to stdout:

```
$ gtts-cli 'hello'
```

Read stdin to hello.mp3 via <text> or <file>:

```
$ echo -n 'hello' | gtts-cli - --output hello.mp3
$ echo -n 'hello' | gtts-cli --file - --output hello.mp3
```

Read 'no check' to nocheck.mp3 without language checking:

```
$ gtts-cli 'no check' --lang zh --nocheck --ouput nocheck.mp3
```

Note: Using `--nocheck` can speed up execution. It exists mostly however to force a <lang> language tag that might not be documented but would work with the API, such as for specific regional sub-tags of documented tags (examples for 'en': 'en-gb', 'en-au', etc.).

1.1.3 Playing sound directly

You can pipe the output of `gtts-cli` into any media player that supports `stdin`. For example, using the `play` command from [SoX](#):

```
$ gtts-cli 'hello' | play -t mp3 -
```

1.2 Module (gtts)

- *gTTS* (`gtts.gTTS`)
- *Languages* (`gtts.lang`)
- *Localized 'accents'*
- *Examples*
- *Playing sound directly*
- *Logging*

1.2.1 gTTS (gtts.gTTS)

```
class gttts.tts.gTTS(text, tld='com', lang='en', slow=False, lang_check=True,  
    pre_processor_funcs=[<function tone_marks>, <function end_of_line>, <function  
abbreviations>, <function word_sub>], tokenizer_func=<bound method Tokenizer.run of  
re.compile('(?=&\n?).|(?<=!).|(?<=).|(?<=).|(?<=\[a-z]\)\ \ |(?<=\[a-z]),  
|(?<=\d):|\.\.\ |\[\[\|\|\];|—|_|!||\|\n\|\|\n\)|' re.IGNORECASE) from: [<function  
tone_marks>, <function period_comma>, <function colon>, <function  
other_punctuation>]>>)
```

gTTS – Google Text-to-Speech.

An interface to Google Translate's Text-to-Speech API.

Parameters

- **text** (*string*) – The text to be read.
- **tld** (*string*) – Top-level domain for the Google Translate host, i.e. `https://translate.google.<tld>`. Different Google domains can produce different localized ‘accents’ for a given language. This is also useful when `google.com` might be blocked within a network but a local or different Google host (e.g. `google.cn`) is not. Default is `com`.
- **lang** (*string, optional*) – The language (IETF language tag) to read the text in. Default is `en`.
- **slow** (*bool, optional*) – Reads text more slowly. Defaults to `False`.
- **lang_check** (*bool, optional*) – Strictly enforce an existing `lang`, to catch a language error early. If set to `True`, a `ValueError` is raised if `lang` doesn’t exist. Setting `lang_check` to `False` skips Web requests (to validate language) and therefore speeds up instantiation. Default is `True`.

- **pre_processor_funcs** (*list*) – A list of zero or more functions that are called to transform (pre-process) text before tokenizing. Those functions must take a string and return a string. Defaults to:

```
[
    pre_processors.tone_marks,
    pre_processors.end_of_line,
    pre_processors.abbreviations,
    pre_processors.word_sub
]
```

- **tokenizer_func** (*callable*) – A function that takes in a string and returns a list of string (tokens). Defaults to:

```
Tokenizer([
    tokenizer_cases.tone_marks,
    tokenizer_cases.period_comma,
    tokenizer_cases.colon,
    tokenizer_cases.other_punctuation
]).run
```

See also:

Pre-processing and tokenizing

Raises

- **AssertionError** – When text is None or empty; when there's nothing left to speak after pre-processing, tokenizing and cleaning.
- **ValueError** – When lang_check is True and lang is not supported.
- **RuntimeError** – When lang_check is True but there's an error loading the languages dictionary.

get_bodies()

Get TTS API request bodies(s) that would be sent to the TTS API.

Returns A list of TTS API request bodies to make.

Return type list

save(*savefile*)

Do the TTS API request and write result to file.

Parameters **savefile** (*string*) – The path and file name to save the mp3 to.

Raises *gTTSError* – When there's an error with the API request.

write_to_fp(*fp*)

Do the TTS API request(s) and write bytes to a file-like object.

Parameters **fp** (*file object*) – Any file-like object to write the mp3 to.

Raises

- *gTTSError* – When there's an error with the API request.
- **TypeError** – When fp is not a file-like object that takes bytes.

exception gttts.tts.gTTSError(*msg=None, **kwargs*)

Exception that uses context to present a meaningful error message

infer_msg(*tts*, *rsp=None*)

Attempt to guess what went wrong by using known information (e.g. http response) and observed behaviour

1.2.2 Languages (gtts.lang)

Note: The easiest way to get a list of available languages is to print them with `gtts-cli --all`

`gtts.lang.tts_langs()`

Languages Google Text-to-Speech supports.

Returns

A dictionary of the type { '<lang>': '<name>' }

Where <lang> is an IETF language tag such as *en* or *zh-TW*, and <name> is the full English name of the language, such as *English* or *Chinese (Mandarin/Taiwan)*.

Return type

The dictionary returned combines languages from two origins:

- Languages fetched from Google Translate (pre-generated in `gtts.langs`)
- Languages that are undocumented variations that were observed to work and present different dialects or accents.

1.2.3 Localized 'accents'

For a given language, Google Translate text-to-speech can speak in different local 'accents' depending on the Google domain (`google.<tld>`) of the request, with some examples shown in the table below.

Note: This is an **incomplete** list. Try different combinaisons of language codes and [known localized Google domains](#). Feel free to add new combinaisons to this list via a Pull Request!

Local accent	Language code (lang)	Top-level domain (tld)
English (Australia)	en	com.au
English (United Kingdom)	en	co.uk
English (United States)	en	com (default)
English (Canada)	en	ca
English (India)	en	co.in
English (Ireland)	en	ie
English (South Africa)	en	co.za
French (Canada)	fr	ca
French (France)	fr	fr
Mandarin (China Mainland)	zh-CN	any
Mandarin (Taiwan)	zh-TW	any
Portuguese (Brazil)	pt	com.br
Portuguese (Portugal)	pt	pt
Spanish (Mexico)	es	com.mx
Spanish (Spain)	es	es
Spanish (United States)	es	com (default)

1.2.4 Examples

Write 'hello' in English to `hello.mp3`:

```
>>> from gtts import gTTS
>>> tts = gTTS('hello', lang='en')
>>> tts.save('hello.mp3')
```

Write 'hello' in Australian English to `hello.mp3`:

```
>>> from gtts import gTTS
>>> tts = gTTS('hello', lang='en', tld='com.au')
>>> tts.save('hello.mp3')
```

Write 'hello bonjour' in English then French to `hello_bonjour.mp3`:

```
>>> from gtts import gTTS
>>> tts_en = gTTS('hello', lang='en')
>>> tts_fr = gTTS('bonjour', lang='fr')
>>>
>>> with open('hello_bonjour.mp3', 'wb') as f:
...     tts_en.write_to_fp(f)
...     tts_fr.write_to_fp(f)
```

1.2.5 Playing sound directly

There's quite a few libraries that do this. Write 'hello' to a file-like object to do further manipulation::

```
>>> from gtts import gTTS
>>> from io import BytesIO
>>>
>>> mp3_fp = BytesIO()
>>> tts = gTTS('hello', lang='en')
>>> tts.write_to_fp(mp3_fp)
>>>
>>> # Load `mp3_fp` as an mp3 file in
>>> # the audio library of your choice
```

Note: See [Issue #26](#) for a discussion and examples of direct playback using various methods.

1.2.6 Logging

gtts does logging using the standard Python logging module. The following loggers are available:

gtts.tts Logger used for the gTTS class

gtts.lang Logger used for the lang module (language fetching)

gtts Upstream logger for all of the above

1.3 Pre-processing and tokenizing

The `gtts.tokenizer` module powers the default pre-processing and tokenizing features of gTTS and provides tools to easily expand them. `gtts.tts.gTTS` takes two arguments `pre_processor_funcs` (list of functions) and `tokenizer_func` (function). See: *Pre-processing, Tokenizing*.

- *Definitions*
- *Pre-processing*
 - *Customizing & Examples*
- *Tokenizing*
 - *Customizing & Examples*
 - *Using a 3rd-party tokenizer*
- *Minimizing*
- *gtts.tokenizer module reference (gtts.tokenizer)*

1.3.1 Definitions

Pre-processor: Function that takes text and returns text. Its goal is to modify text (for example correcting pronunciation), and/or to prepare text for proper tokenization (for example ensuring spacing after certain characters).

Tokenizer: Function that takes text and returns it split into a list of *tokens* (strings). In the gTTS context, its goal is to cut the text into smaller segments that do not exceed the maximum character size allowed for each TTS API request, while making the speech sound natural and continuous. It does so by splitting text where speech would naturally pause (for example on “.”) while handling where it should not (for example on “10.5” or “U.S.A.”). Such rules are called *tokenizer cases*, which it takes a list of.

Tokenizer case: Function that defines one of the specific cases used by `gtts.tokenizer.core.Tokenizer`. More specifically, it returns a `regex` object that describes what to look for for a particular case. `gtts.tokenizer.core.Tokenizer` then creates its main *regex* pattern by joining all *tokenizer cases* with “[”.

1.3.2 Pre-processing

You can pass a list of any function to `gtts.tts.gTTS`’s `pre_processor_funcs` attribute to act as pre-processor (as long as it takes a string and returns a string).

By default, `gtts.tts.gTTS` takes a list of the following pre-processors, applied in order:

```
[
    pre_processors.tone_marks,
    pre_processors.end_of_line,
    pre_processors.abbreviations,
    pre_processors.word_sub
]
```

`gtts.tokenizer.pre_processors.abbreviations(text)`

Remove periods after an abbreviation from a list of known abbreviations that can be spoken the same without that period. This prevents having to handle tokenization of that period.

Note: Could potentially remove the ending period of a sentence.

Note: Abbreviations that Google Translate can't pronounce without (or even with) a period should be added as a word substitution with a `PreProcessorSub` pre-processor. Ex.: 'Esq.', 'Esquire'.

`gtts.tokenizer.pre_processors.end_of_line(text)`
Re-form words cut by end-of-line hyphens.

Remove "<hyphen><newline>".

`gtts.tokenizer.pre_processors.tone_marks(text)`
Add a space after tone-modifying punctuation.

Because the `tone_marks` tokenizer case will split after a tone-modifying punctuation mark, make sure there's whitespace after.

`gtts.tokenizer.pre_processors.word_sub(text)`
Word-for-word substitutions.

Customizing & Examples

This module provides two classes to help build pre-processors:

- `gtts.tokenizer.core.PreProcessorRegex` (for *regex*-based replacing, as would `re.sub` use)
- `gtts.tokenizer.core.PreProcessorSub` (for word-for-word replacements).

The `run(text)` method of those objects returns the processed text.

Speech corrections (word substitution)

The default substitutions are defined by the `gtts.tokenizer.symbols.SUB_PAIRS` list. Add a custom one by appending to it:

```
>>> from gtts.tokenizer import pre_processors
>>> import gtts.tokenizer.symbols
>>>
>>> gtts.tokenizer.symbols.SUB_PAIRS.append(
...     ('sub.', 'submarine')
... )
>>> test_text = "Have you seen the Queen's new sub.?"
>>> pre_processors.word_sub(test_text)
"Have you seen the Queen's new submarine?"
```

Abbreviations

The default abbreviations are defined by the `gtts.tokenizer.symbols.ABBREVIATIONS` list. Add a custom one to it to add a new abbreviation to remove the period from. *Note: the default list already includes an extensive list of English abbreviations that Google Translate will read even without the period.*

See `gtts.tokenizer.pre_processors` for more examples.

1.3.3 Tokenizing

You can pass any function to `gtts.tts.gTTS`'s `tokenizer_func` attribute to act as tokenizer (as long as it takes a string and returns a list of strings).

By default, gTTS takes the `gtts.tokenizer.core.Tokenizer`'s `gtts.tokenizer.core.Tokenizer.run()`, initialized with default *tokenizer cases*:

```
Tokenizer([
    tokenizer_cases.tone_marks,
    tokenizer_cases.period_comma,
    tokenizer_cases.other_punctuation
]).run
```

The available *tokenizer cases* are as follows:

`gtts.tokenizer.tokenizer_cases.colon()`

Colon case.

Match a colon “:” only if not preceded by a digit. Mainly to prevent a cut in the middle of time notations e.g. 10:01

`gtts.tokenizer.tokenizer_cases.legacy_all_punctuation()`

Match all punctuation.

Use as only tokenizer case to mimic gTTS 1.x tokenization.

`gtts.tokenizer.tokenizer_cases.other_punctuation()`

Match other punctuation.

Match other punctuation to split on; punctuation that naturally inserts a break in speech.

`gtts.tokenizer.tokenizer_cases.period_comma()`

Period and comma case.

Match if not preceded by “.<letter>” and only if followed by space. Won't cut in the middle/after dotted abbreviations; won't cut numbers.

Note: Won't match if a dotted abbreviation ends a sentence.

Note: Won't match the end of a sentence if not followed by a space.

`gtts.tokenizer.tokenizer_cases.tone_marks()`

Keep tone-modifying punctuation by matching following character.

Assumes the *tone_marks* pre-processor was run for cases where there might not be any space after a tone-modifying punctuation mark.

Customizing & Examples

A *tokenizer case* is a function that returns a compiled *regex* object to be used in a `re.split()` context.

`gtts.tokenizer.core.Tokenizer` takes a list of *tokenizer cases* and joins their pattern with “|” in one single pattern.

This module provides a class to help build tokenizer cases: `gtts.tokenizer.core.RegexBuilder`. See `gtts.tokenizer.core.RegexBuilder` and `gtts.tokenizer.tokenizer_cases` for examples.

Using a 3rd-party tokenizer

Even though `gtts.tokenizer.core.Tokenizer` works well in this context, there are way more advanced tokenizers and tokenizing techniques. As long as you can restrict the length of output tokens, you can use any tokenizer you’d like, such as the ones in [NLTK](#).

1.3.4 Minimizing

The Google Translate text-to-speech API accepts a maximum of **100 characters**.

If after tokenization any of the tokens is larger than 100 characters, it will be split in two:

- On the last space character that is closest to, but before the 100th character;
- Between the 100th and 101st characters if there’s no space.

1.3.5 gtts.tokenizer module reference (gtts.tokenizer)

class `gtts.tokenizer.core.RegexBuilder`(*pattern_args*, *pattern_func*, *flags=0*)

Builds regex using arguments passed into a pattern template.

Builds a regex object for which the pattern is made from an argument passed into a template. If more than one argument is passed (iterable), each pattern is joined by “|” (regex alternation ‘or’) to create a single pattern.

Parameters

- **pattern_args** (*iteratable*) – String element(s) to be each passed to `pattern_func` to create a regex pattern. Each element is `re.escape`’d before being passed.
- **pattern_func** (*callable*) – A ‘template’ function that should take a string and return a string. It should take an element of `pattern_args` and return a valid regex pattern group string.
- **flags** – `re` flag(s) to compile with the regex.

Example

To create a simple regex that matches on the characters “a”, “b”, or “c”, followed by a period:

```
>>> rb = RegexBuilder('abc', lambda x: "{}\.".format(x))
```

Looking at `rb.regex` we get the following compiled regex:

```
>>> print(rb.regex)
'a\.|b\.|c\.'
```


The above is fairly simple, but this class can help in writing more complex repetitive regex, making them more readable and easier to create by using existing data structures.

Example

To match the character following the words “lorem”, “ipsum”, “meili” or “koda”:

```
>>> words = ['lorem', 'ipsum', 'meili', 'koda']
>>> rb = RegexBuilder(words, lambda x: "(?<={}).".format(x))
```

Looking at `rb.regex` we get the following compiled regex:

```
>>> print(rb.regex)
'(?<=lorem).|(?<=ipsum).|(?<=meili).|(?<=koda).'
```

class `gtts.tokenizer.core.PreProcessorRegex`(*search_args*, *search_func*, *repl*, *flags=0*)

Regex-based substitution text pre-processor.

Runs a series of regex substitutions (`re.sub`) from each `regex` of a `gtts.tokenizer.core.RegexBuilder` with an extra `repl` replacement parameter.

Parameters

- **search_args** (*iteratable*) – String element(s) to be each passed to `search_func` to create a regex pattern. Each element is `re.escape`’d before being passed.
- **search_func** (*callable*) – A ‘template’ function that should take a string and return a string. It should take an element of `search_args` and return a valid regex search pattern string.
- **repl** (*string*) – The common replacement passed to the `sub` method for each `regex`. Can be a raw string (the case of a regex backreference, for example)
- **flags** – `re` flag(s) to compile with each *regex*.

Example

Add “!” after the words “lorem” or “ipsum”, while ignoring case:

```
>>> import re
>>> words = ['lorem', 'ipsum']
>>> pp = PreProcessorRegex(words,
...                         lambda x: "({})".format(x), r'\!1!',
...                         re.IGNORECASE)
```

In this case, the `regex` is a group and the replacement uses its backreference `\1` (as a raw string). Looking at `pp` we get the following list of search/replacement pairs:

```
>>> print(pp)
(re.compile('(lorem)', re.IGNORECASE), repl='\!1'),
(re.compile('(ipsum)', re.IGNORECASE), repl='\!1')
```

It can then be run on any string of text:

```
>>> pp.run("LOREM ipSuM")
"LOREM! ipSuM!"
```

See [gtts.tokenizer.pre_processors](#) for more examples.

run(text)

Run each regex substitution on `text`.

Parameters `text (string)` – the input text.

Returns text after all substitutions have been sequentially applied.

Return type string

class `gtts.tokenizer.core.PreProcessorSub(sub_pairs, ignore_case=True)`

Simple substitution text preprocessor.

Performs string-for-string substitution from list a find/replace pairs. It abstracts [gtts.tokenizer.core.PreProcessorRegex](#) with a default simple substitution regex.

Parameters

- **sub_pairs (list)** – A list of tuples of the style (<search str>, <replace str>)
- **ignore_case (bool)** – Ignore case during search. Defaults to True.

Example

Replace all occurrences of “Mac” to “PC” and “Firefox” to “Chrome”:

```
>>> sub_pairs = [('Mac', 'PC'), ('Firefox', 'Chrome')]
>>> pp = PreProcessorSub(sub_pairs)
```

Looking at the `pp`, we get the following list of search (regex)/replacement pairs:

```
>>> print(pp)
(re.compile('Mac', re.IGNORECASE), repl='PC'),
(re.compile('Firefox', re.IGNORECASE), repl='Chrome')
```

It can then be run on any string of text:

```
>>> pp.run("I use firefox on my mac")
"I use Chrome on my PC"
```

See [gtts.tokenizer.pre_processors](#) for more examples.

run(text)

Run each substitution on `text`.

Parameters `text (string)` – the input text.

Returns text after all substitutions have been sequentially applied.

Return type string

class `gtts.tokenizer.core.Tokenizer(regex_funcs, flags=<RegexFlag.IGNORECASE: 2>)`

An extensible but simple generic rule-based tokenizer.

A generic and simple string tokenizer that takes a list of functions (called *tokenizer cases*) returning `regex` objects and joins them by “|” (regex alternation ‘or’) to create a single regex to use with the standard `regex.split()` function.

`regex_funcs` is a list of any function that can return a `regex` (from `re.compile()`) object, such as a [gtts.tokenizer.core.RegexBuilder](#) instance (and its `regex` attribute).

See the `gtts.tokenizer.tokenizer_cases` module for examples.

Parameters

- **regex_funcs** (*list*) – List of compiled `regex` objects. Each functions's pattern will be joined into a single pattern and compiled.
- **flags** – `re` flag(s) to compile with the final `regex`. Defaults to `re.IGNORECASE`

Note: When the `regex` objects obtained from `regex_funcs` are joined, their individual `re` flags are ignored in favour of `flags`.

Raises `TypeError` – When an element of `regex_funcs` is not a function, or a function that does not return a compiled `regex` object.

Warning: Joined `regex` patterns can easily interfere with one another in unexpected ways. It is recommended that each tokenizer case operate on distinct or non-overlapping characters/sets of characters (For example, a tokenizer case for the period (“.”) should also handle not matching/cutting on decimals, instead of making that a separate tokenizer case).

Example

A tokenizer with a two simple case (*Note: these are bad cases to tokenize on, this is simply a usage example*):

```
>>> import re, RegexBuilder
>>>
>>> def case1():
...     return re.compile("\\,")
>>>
>>> def case2():
...     return RegexBuilder('abc', lambda x: "{}\\.".format(x)).regex
>>>
>>> t = Tokenizer([case1, case2])
```

Looking at `case1().pattern`, we get:

```
>>> print(case1().pattern)
'\\,'
```

Looking at `case2().pattern`, we get:

```
>>> print(case2().pattern)
'a\\.|b\\.|c\\.'
```

Finally, looking at `t`, we get them combined:

```
>>> print(t)
're.compile('\\,|a\\.|b\\.|c\\.', re.IGNORECASE)
from: [<function case1 at 0x10bbcd08>, <function case2 at 0x10b5c5e18>']
```

It can then be run on any string of text:

```
>>> t.run("Hello, my name is Linda a. Call me Lin, b. I'm your friend")
['Hello', ' my name is Linda ', ' Call me Lin', ' ', " I'm your friend"]
```

run(*text*)

Tokenize *text*.

Parameters **text** (*string*) – the input text to tokenize.

Returns A list of strings (token) split according to the tokenizer cases.

Return type list

```
symbols.ABBREVIATIONS = ['dr', 'jr', 'mr', 'mrs', 'ms', 'msgr', 'prof', 'sr', 'st']
```

```
symbols.SUB_PAIRS = [('Esq.', 'Esquire')]
```

```
symbols.ALL_PUNC = '?!.,;()[]<...;-\\n'
```

```
symbols.TONE_MARKS = '?!'
```

1.4 License

The MIT License (MIT)

Copyright © 2014-2021 Pierre Nicolas Durette

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the “Software”), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED “AS IS”, WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

1.5 Contributing

1.5.1 Reporting Issues

On the Github [issues](#) page. Thanks!

1.5.2 Submitting Patches

1. **Fork.** Follow [PEP 8](#)!
2. **Write/Update tests** (see below).
3. **Document.** Docstrings follow the [Google Python Style Guide](#) (docs by [Sphinx](#)). You can ‘test’ documentation:

```
$ pip install .[docs]
$ cd docs && make html # generated in docs/_build/html/
```

4. **Open Pull Request.** To the master branch.
5. **Changelog.** This project uses [towncrier](#) for managing the changelog. Please consider creating one or more ‘news fragment’ in the `/news/` directory and adding them to your PR, in the style of `<issue_or_pr_number>.<type>` where ‘type’ is one of: ‘feature’, ‘bugfix’, ‘doc’, ‘removal’ or ‘misc’.

See [towncrier](#) (New Fragments) for more details. Example:

```
$ echo 'Fixed a thing!' > gtts/news/1234.bugfix
```

Note:

Please don’t hesitate to contribute! While good tests, docs and structure are encouraged, I do welcome great ideas over absolute conformity to the above!

Thanks!

1.5.3 Testing

Testing is done with the `unittest` framework.

As a rule, the file `./tests/test_<module>.py` file tests the `<module>` module.

To run all tests (testing only language ‘en’ and generating an html coverage report in `gtts/htmlcov/`):

```
$ pip install .[tests]
$ TEST_LANGS=en pytest -v -s gtts/ --cov=gtts --cov-report=html
```

1.6 Changelog

1.6.1 2.2.3 (2021-06-17)

Features

- Added Bulgarian language support ([#302](#))

1.6.2 2.2.2 (2021-02-03)

Features

- Adds a language fallback feature for deprecated languages to maintain compatibility (e.g. `en-us` becomes `en`). Fallback can be disabled with `lang_check=False` or `--nocheck` for the cli (#267)

Bugfixes

- Fix Python 2.7 compatibility (!). Python 2 is long gone, but the cut wasn't clearly communicated for gTTS, so it was restored. Python 2 support will be completely removed in the next major release. (#255)
- Language code case sensitivity is maintained throughout (#267)

Deprecations and Removals

- The following list of 'hyphenated' language codes no longer work and have been removed: `en-us`, `en-ca`, `en-uk`, `en-gb`, `en-au`, `en-gh`, `en-in`, `en-ie`, `en-nz`, `en-ng`, `en-ph`, `en-za`, `en-tz`, `fr-ca`, `fr-fr`, `pt-br`, `pt-pt`, `es-es`, `es-us`, `zh-cn`, `zh-tw` (#267)
- Removed the `gtts.get_url()` method (outdated since 2.1.0) (#270)

1.6.3 2.2.1 (2020-11-15)

Bugfixes

- `_package_rpc()` was erroneously packaging the entire text instead of tokenized part (#252)

Improved Documentation

- Removes reference to automatic retrieval of languages (#250)

Misc

- #251

1.6.4 2.2.0 (2020-11-14)

Features

- Switch to the newer Google TTS API (thanks to @Boudewijn26!). See [his great writeup](#) for more on the methodology and why this was necessary. (#226, #232, #236, #241)

Deprecations and Removals

- Removed automatic language download from the main code, which has become too unreliable & slow. Languages will still be fetched but a pre-generated list will be shipped with gTTS. (#233, #241, #242, #243)
- Because languages are now pre-generated, removed custom TLD support for language URL (which allowed to get language **names** in other than English) (#245)

Misc

- #245

1.6.5 2.1.2 (2020-11-10)

Features

- Update *gTTS-token* to 1.1.4 (#238)

Bugfixes

- Fixed an issue where some tokens could be empty after minimization (#229, #239)

Improved Documentation

- Grammar, spelling and example fixes (#227)

Misc

- #218, #230, #231, #239

1.6.6 2.1.1 (2020-01-25)

Bugfixes

- Debug mode now uses a copy of `locals()` to prevent `RuntimeError` (#213)

1.6.7 2.1.0 (2020-01-01)

Features

- The `gtts` module
 - Added the ability to customize the Google Translate URL hostname. This is useful when `google.com` might be blocked within a network but a local or different Google host (e.g. `google.cn`) is not (#143, #203):
 - * New `gTTS()` parameter `tld` to specify the top-level domain to use for the Google hostname, i.e. `https://translate.google.<tld>` (default: `com`).
 - * Languages are also now fetched using the same customized hostname.

- Pre-generated TTS API request URLs can now be obtained instead of writing an `mp3` file to disk (for example to be used in an external program):
 - * New `get_urls()` method returns the list of URLs generated by gTTS, which can be used in lieu of `write_to_fp()` or `save()`.
- The `gtts-cli` command-line tool
 - New `--tld` option to match the new `gtts` customizable hostname (#200, #207)
- Other
 - Added Python 3.8 support (#204)

Bugfixes

- Changed default word-for-word pre-processor (`('M.', 'Monsieur')`) which would substitute any `'m.'` for `'monsieur'` (e.g. `'them.'` became `'themonsieur'`) (#197)

Improved Documentation

- Added examples for newer features (#205, #207)

Misc

- #204, #205, #207

1.6.8 2.0.4 (2019-08-29)

Features

- gTTS is now built as a wheel package (Python 2 & 3) (#181)

Improved Documentation

- Fixed bad example in docs (#163, #166)

Misc

- #164, #171, #173, #185

1.6.9 2.0.3 (2018-12-15)

Features

- Added new tokenizer case for `':'` preventing cut in the middle of a time notation (#135)

Misc

- [#159](#)

1.6.10 2.0.2 (2018-12-09)

Features

- Added Python 3.7 support, modernization of packaging, testing and CI ([#126](#))

Bugfixes

- Fixed language retrieval/validation broken from new Google Translate page ([#156](#))

1.6.11 2.0.1 (2018-06-20)

Bugfixes

- Fixed an UnicodeDecodeError when installing gTTS if system locale was not utf-8 ([#120](#))

Improved Documentation

- Added *Pre-processing and tokenizing > Minimizing* section about the API's 100 characters limit and how larger tokens are handled ([#121](#))

Misc

- [#122](#)

1.6.12 2.0.0 (2018-04-30)

([#108](#))

Features

- The `gtts` module
 - New logger ("`gtts`") replaces all occurrences of `print()`
 - Languages list is now obtained automatically (`gtts.lang`) ([#91](#), [#94](#), [#106](#))
 - Added a curated list of language sub-tags that have been observed to provide different dialects or accents (e.g. "`en-gb`", "`fr-ca`")
 - New `gTTS()` parameter `lang_check` to disable language checking.
 - `gTTS()` now delegates the `text` tokenizing to the API request methods (i.e. `write_to_fp()`, `save()`), allowing `gTTS` instances to be modified/reused
 - Rewrote tokenizing and added pre-processing (see below)

- New `gTTS()` parameters `pre_processor_funcs` and `tokenizer_func` to configure pre-processing and tokenizing (or use a 3rd party tokenizer)
- Error handling:
 - * Added new exception `gTTSError` raised on API request errors. It attempts to guess what went wrong based on known information and observed behaviour (#60, #106)
 - * `gTTS.write_to_fp()` and `gTTS.save()` also raise `gTTSError` on *gtts_token* error
 - * `gTTS.write_to_fp()` raises `TypeError` when `fp` is not a file-like object or one that doesn't take bytes
 - * `gTTS()` raises `ValueError` on unsupported languages (and `lang_check` is `True`)
 - * More fine-grained error handling throughout (e.g. *request failed* vs. *request successful with a bad response*)
- Tokenizer (and new pre-processors):
 - Rewrote and greatly expanded tokenizer (`gtts.tokenizer`)
 - Smarter token 'cleaning' that will remove tokens that only contain characters that can't be spoken (i.e. punctuation and whitespace)
 - Decoupled token minimizing from tokenizing, making the latter usable in other contexts
 - New flexible speech-centric text pre-processing
 - New flexible full-featured regex-based tokenizer (`gtts.tokenizer.core.Tokenizer`)
 - New `RegexBuilder`, `PreProcessorRegex` and `PreProcessorSub` classes to make writing regex-powered text *pre-processors* and *tokenizer cases* easier
 - Pre-processors:
 - * Re-form words cut by end-of-line hyphens
 - * Remove periods after a (customizable) list of known abbreviations (e.g. "jr", "sr", "dr") that can be spoken the same without a period
 - * Perform speech corrections by doing word-for-word replacements from a (customizable) list of tuples
 - Tokenizing:
 - * Keep punctuation that modify the inflection of speech (e.g. "?", "!")
 - * Don't split in the middle of numbers (e.g. "10.5", "20,000,000") (#101)
 - * Don't split on "dotted" abbreviations and acronyms (e.g. "U.S.A")
 - * Added Chinese comma ("，"), ellipsis ("…") to punctuation list to tokenize on (#86)
- The `gtts-cli` command-line tool
 - Rewrote cli as first-class citizen module (`gtts.cli`), powered by [Click](#)
 - Windows support using *setuptools's entry_points*
 - Better support for Unicode I/O in Python 2
 - All arguments are now pre-validated
 - New `--nocheck` flag to skip language pre-checking
 - New `--all` flag to list all available languages
 - Either the `--file` option or the `<text>` argument can be set to `"-"` to read from `stdin`

- The `--debug` flag uses logging and doesn't pollute `stdout` anymore

Bugfixes

- `_minimize()`: Fixed an infinite recursion loop that would occur when a token started with the minimizing delimiter (i.e. a space) (#86)
- `_minimize()`: Handle the case where a token of more than 100 characters did not contain a space (e.g. in Chinese).
- Fixed an issue that fused multiline text together if the total number of characters was less than 100
- Fixed `gtts-cli` Unicode errors in Python 2.7 (famous last words) (#78, #93, #96)

Deprecations and Removals

- Dropped Python 3.3 support
- Removed `debug` parameter of `gTTS` (in favour of `logger`)
- `gtts-cli`: Changed long option name of `-o` to `--output` instead of `--destination`
- `gTTS()` will raise a `ValueError` rather than an `AssertionError` on unsupported language

Improved Documentation

- Rewrote all documentation files as `reStructuredText`
- Comprehensive documentation written for `Sphinx`, published to <http://gtts.readthedocs.io>
- Changelog built with `towncrier`

Misc

- Major test re-work
- Language tests can read a `TEST_LANGS` environment variable so not all language tests are run every time.
- Added `AppVeyor` CI for Windows
- `PEP 8` compliance

1.6.13 1.2.2 (2017-08-15)

Misc

- Update LICENCE, add to manifest (#77)

1.6.14 1.2.1 (2017-08-02)

Features

- Add Unicode punctuation to the tokenizer (such as for Chinese and Japanese) (#75)

Bugfixes

- Fix > 100 characters non-ASCII split, `unicode()` for Python 2 (#71, #73, #75)

1.6.15 1.2.0 (2017-04-15)

Features

- Option for slower read speed (`slow=True` for `gTTS()`, `--slow` for `gtts-cli`) (#40, #41, #64, #67)
- System proxy settings are passed transparently to all http requests (#45, #68)
- Silence SSL warnings from `urllib3` (#69)

Bugfixes

- The text to read is now cut in proper chunks in Python 2 unicode. This broke reading for many languages such as Russian.
- Disabled SSL verify on http requests to accommodate certain firewalls and proxies.
- Better Python 2/3 support in general (#9, #48, #68)

Deprecations and Removals

- `'pt-br'` : 'Portuguese (Brazil)' (it was the same as `'pt'` and not Brazilian) (#69)

1.6.16 1.1.8 (2017-01-15)

Features

- Added stdin support via the `'-'` text argument to `gtts-cli` (#56)

1.6.17 1.1.7 (2016-12-14)

Features

- Added utf-8 support to `gtts-cli` (#52)

1.6.18 1.1.6 (2016-07-20)

Features

- Added 'bn' : 'Bengali' (#39, #44)

Deprecations and Removals

- 'ht' : 'Haitian Creole' (removed by Google) (#43)

1.6.19 1.1.5 (2016-05-13)

Bugfixes

- Fixed HTTP 403s by updating the client argument to reflect new API usage (#32, #33)

1.6.20 1.1.4 (2016-02-22)

Features

- Spun-off token calculation to gTTS-Token (#23, #29)

1.6.21 1.1.3 (2016-01-24)

Bugfixes

- gtts-cli works with Python 3 (#20)
- Better support for non-ASCII characters (#21, #22)

Misc

- Moved out gTTS token to its own module (#19)

1.6.22 1.1.2 (2016-01-13)

Features

- Added gTTS token (tk url parameter) calculation (#14, #15, #17)

1.6.23 1.0.7 (2015-10-07)

Features

- Added stdout support to `gtts-cli`, text now an argument rather than an option (#10)

1.6.24 1.0.6 (2015-07-30)

Features

- Raise an exception on bad HTTP response (4xx or 5xx) (#8)

Bugfixes

- Added `client=t` parameter for the api HTTP request (#8)

1.6.25 1.0.5 (2015-07-15)

Features

- `write_to_fp()` to write to a file-like object (#6)

1.6.26 1.0.4 (2015-05-11)

Features

- Added Languages: *zh-yue* : ‘Chinese (Cantonese)’, *en-uk* : ‘English (United Kingdom)’, *pt-br* : ‘Portuguese (Brazil)’, *es-es* : ‘Spanish (Spain)’, *es-us* : ‘Spanish (United StateS)’, *zh-cn* : ‘Chinese (Mandarin/China)’, *zh-tw* : ‘Chinese (Mandarin/Taiwan)’ (#4)

Bugfixes

- `gtts-cli` print version and pretty printed available languages, language codes are now case insensitive (#4)

1.6.27 1.0.3 (2014-11-21)

Features

- Added Languages: ‘en-us’ : ‘English (United States)’, ‘en-au’ : ‘English (Australia)’ (#3)

1.6.28 1.0.2 (2014-05-15)

Features

- Python 3 support

1.6.29 1.0.1 (2014-05-15)

Misc

- SemVer versioning, CI changes

1.6.30 1.0 (2014-05-08)

Features

- Initial release

CHAPTER
TWO

MISC

- genindex
- modindex

PYTHON MODULE INDEX

g

- `gtts.lang`, [7](#)
- `gtts.tokenizer`, [8](#)
- `gtts.tokenizer.pre_processors`, [9](#)
- `gtts.tokenizer.tokenizer_cases`, [11](#)
- `gtts.tts`, [5](#)

Symbols

--all
 gtts-cli command line option, 3
 --debug
 gtts-cli command line option, 4
 --file <file>
 gtts-cli command line option, 3
 --lang <lang>
 gtts-cli command line option, 3
 --nocheck
 gtts-cli command line option, 3
 --output <file>
 gtts-cli command line option, 3
 --slow
 gtts-cli command line option, 3
 --tld <tld>
 gtts-cli command line option, 3
 --version
 gtts-cli command line option, 4
 -f
 gtts-cli command line option, 3
 -l
 gtts-cli command line option, 3
 -o
 gtts-cli command line option, 3
 -s
 gtts-cli command line option, 3
 -t
 gtts-cli command line option, 3
 <text>
 gtts-cli command line option, 4

A

ABBREVIATIONS (*gtts.tokenizer.symbols attribute*), 16
 abbreviations() (in *module gtts.tokenizer.pre_processors*), 9
 ALL_PUNC (*gtts.tokenizer.symbols attribute*), 16

C

colon() (in *module gtts.tokenizer.tokenizer_cases*), 11

E

end_of_line() (in *module gtts.tokenizer.pre_processors*), 10

G

get_bodies() (*gtts.tts.gTTS method*), 6
 gTTS (*class in gtts.tts*), 5
 gtts.lang
 module, 7
 gtts.tokenizer
 module, 8
 gtts.tokenizer.pre_processors
 module, 9
 gtts.tokenizer.tokenizer_cases
 module, 11
 gtts.tts
 module, 5
 gtts-cli command line option
 --all, 3
 --debug, 4
 --file <file>, 3
 --lang <lang>, 3
 --nocheck, 3
 --output <file>, 3
 --slow, 3
 --tld <tld>, 3
 --version, 4
 -f, 3
 -l, 3
 -o, 3
 -s, 3
 -t, 3
 <text>, 4
 gTTSError, 6

I

infer_msg() (*gtts.tts.gTTSError method*), 6

L

legacy_all_punctuation() (in *module gtts.tokenizer.tokenizer_cases*), 11

M

module

- `gtts.lang`, 7
- `gtts.tokenizer`, 8
- `gtts.tokenizer.pre_processors`, 9
- `gtts.tokenizer.tokenizer_cases`, 11
- `gtts.tts`, 5

O

`other_punctuation()` (in module `gtts.tokenizer.tokenizer_cases`), 11

P

`period_comma()` (in module `gtts.tokenizer.tokenizer_cases`), 11

`PreProcessorRegex` (class in `gtts.tokenizer.core`), 13

`PreProcessorSub` (class in `gtts.tokenizer.core`), 14

R

`RegexBuilder` (class in `gtts.tokenizer.core`), 12

`run()` (`gtts.tokenizer.core.PreProcessorRegex` method), 14

`run()` (`gtts.tokenizer.core.PreProcessorSub` method), 14

`run()` (`gtts.tokenizer.core.Tokenizer` method), 16

S

`save()` (`gtts.tts.gTTS` method), 6

`SUB_PAIRS` (`gtts.tokenizer.symbols` attribute), 16

T

`Tokenizer` (class in `gtts.tokenizer.core`), 14

`TONE_MARKS` (`gtts.tokenizer.symbols` attribute), 16

`tone_marks()` (in module `gtts.tokenizer.pre_processors`), 10

`tone_marks()` (in module `gtts.tokenizer.tokenizer_cases`), 11

`tts_langs()` (in module `gtts.lang`), 7

W

`word_sub()` (in module `gtts.tokenizer.pre_processors`), 10

`write_to_fp()` (`gtts.tts.gTTS` method), 6