

# 1 安装

```
pip install flashtext
```

# 2 提取关键字

```
1 >>> from flashtext import KeywordProcessor
2 >>> keyword_processor = KeywordProcessor()
3 >>> # keyword_processor.add_keyword(<unclean name>, <standardised name>)
4 >>> keyword_processor.add_keyword('Big Apple', 'New York')
5 >>> keyword_processor.add_keyword('Bay Area')
6 >>> keywords_found = keyword_processor.extract_keywords('I love Big Apple and Bay
↳ Area.')
7 >>> keywords_found
8 >>> # ['New York']
```

# 3 替换关键字

```
1 >>> keyword_processor.add_keyword('New Delhi', 'NCR region')
2 >>> new_sentence = keyword_processor.replace_keywords('I love Big Apple and new
↳ delhi.')
3 >>> new_sentence
4 >>> # 'I love New York and NCR region.'
```

## 3.1 Case Sensitive example

```
1 >>> from flashtext import KeywordProcessor
2 >>> keyword_processor = KeywordProcessor(case_sensitive=True)
3 >>> keyword_processor.add_keyword('Big Apple', 'New York')
4 >>> keyword_processor.add_keyword('Bay Area')
5 >>> keywords_found = keyword_processor.extract_keywords('I love big Apple and Bay
↳ Area.')
6 >>> keywords_found
7 >>> # ['Bay Area']
```

## 3.2 关键字空格信息

```
1 >>> from flashtext import KeywordProcessor
2 >>> keyword_processor = KeywordProcessor()
3 >>> keyword_processor.add_keyword('Big Apple', 'New York')
4 >>> keyword_processor.add_keyword('Bay Area')
5 >>> keywords_found = keyword_processor.extract_keywords('I love big Apple and Bay
  ↳ Area.', span_info=True)
6 >>> keywords_found
7 >>> # [('New York', 7, 16), ('Bay Area', 21, 29)]
```

## 3.3 用关键字提取额外的信息

```
1 >>> from flashtext import KeywordProcessor
2 >>> kp = KeywordProcessor()
3 >>> kp.add_keyword('Taj Mahal', ('Monument', 'Taj Mahal'))
4 >>> kp.add_keyword('Delhi', ('Location', 'Delhi'))
5 >>> kp.extract_keywords('Taj Mahal is in Delhi.')
6 >>> # [('Monument', 'Taj Mahal'), ('Location', 'Delhi')]
7 >>> # NOTE: replace_keywords feature won't work with this.
```

No clean name for Keywords

```
1 >>> from flashtext import KeywordProcessor
2 >>> keyword_processor = KeywordProcessor()
3 >>> keyword_processor.add_keyword('Big Apple')
4 >>> keyword_processor.add_keyword('Bay Area')
5 >>> keywords_found = keyword_processor.extract_keywords('I love big Apple and Bay
  ↳ Area.')
6 >>> keywords_found
7 >>> # ['Big Apple', 'Bay Area']
```

## 3.4 同时添加多个关键字

```
1 >>> from flashtext import KeywordProcessor
2 >>> keyword_processor = KeywordProcessor()
3 >>> keyword_dict = {
```

```

4 >>>     "java": ["java_2e", "java programing"],
5 >>>     "product management": ["PM", "product manager"]
6 >>> }
7 >>> # {'clean_name': ['list of unclean names']}
8 >>> keyword_processor.add_keywords_from_dict(keyword_dict)
9 >>> # Or add keywords from a list:
10 >>> keyword_processor.add_keywords_from_list(["java", "python"])
11 >>> keyword_processor.extract_keywords('I am a product manager for a java_2e
    ↪ platform')
12 >>> # output ['product management', 'java']

```

### 3.5 移除关键字

```

1 >>> from flashtext import KeywordProcessor
2 >>> keyword_processor = KeywordProcessor()
3 >>> keyword_dict = {
4 >>>     "java": ["java_2e", "java programing"],
5 >>>     "product management": ["PM", "product manager"]
6 >>> }
7 >>> keyword_processor.add_keywords_from_dict(keyword_dict)
8 >>> print(keyword_processor.extract_keywords('I am a product manager for a java_2e
    ↪ platform'))
9 >>> # output ['product management', 'java']
10 >>> keyword_processor.remove_keyword('java_2e')
11 >>> # you can also remove keywords from a list/ dictionary
12 >>> keyword_processor.remove_keywords_from_dict({"product management": ["PM"]})
13 >>> keyword_processor.remove_keywords_from_list(["java programing"])
14 >>> keyword_processor.extract_keywords('I am a product manager for a java_2e
    ↪ platform')
15 >>> # output ['product management']

```

### 3.6 To check Number of terms in KeywordProcessor

```

1 >>> from flashtext import KeywordProcessor
2 >>> keyword_processor = KeywordProcessor()
3 >>> keyword_dict = {
4 >>>     "java": ["java_2e", "java programing"],

```

```

5 >>> "product management": ["PM", "product manager"]
6 >>> }
7 >>> keyword_processor.add_keywords_from_dict(keyword_dict)
8 >>> print(len(keyword_processor))
9 >>> # output 4

```

### 3.7 To check if term is present in KeywordProcessor

```

1 >>> from flashtext import KeywordProcessor
2 >>> keyword_processor = KeywordProcessor()
3 >>> keyword_processor.add_keyword('j2ee', 'Java')
4 >>> 'j2ee' in keyword_processor
5 >>> # output: True
6 >>> keyword_processor.get_keyword('j2ee')
7 >>> # output: Java
8 >>> keyword_processor['colour'] = 'color'
9 >>> keyword_processor['colour']
10 >>> # output: color

```

### 3.8 Get all keywords in dictionary

```

1 >>> from flashtext import KeywordProcessor
2 >>> keyword_processor = KeywordProcessor()
3 >>> keyword_processor.add_keyword('j2ee', 'Java')
4 >>> keyword_processor.add_keyword('colour', 'color')
5 >>> keyword_processor.get_all_keywords()
6 >>> # output: {'colour': 'color', 'j2ee': 'Java'}

```

For detecting Word Boundary currently any character other than this is considered a word boundary.

To set or add characters as part of word characters

```

1 >>> from flashtext import KeywordProcessor
2 >>> keyword_processor = KeywordProcessor()
3 >>> keyword_processor.add_keyword('Big Apple')
4 >>> print(keyword_processor.extract_keywords('I love Big Apple/Bay Area.'))
5 >>> # ['Big Apple']

```

```

6 >>> keyword_processor.add_non_word_boundary('/')
7 >>> print(keyword_processor.extract_keywords('I love Big Apple/Bay Area.'))
8 >>> # []

```

API doc KeywordProcessor

API Doc Import and initialize module Add Keywords to module Extract keywords Replace keywords Add keywords from File Add keywords from dict Add keywords from list KeywordProcessor Class Doc Test

```

1 $ git clone https://github.com/vi3k6i5/flashtext
2 $ cd flashtext
3 $ pip install pytest
4 $ python setup.py test

```

Build Docs

```

1 $ git clone https://github.com/vi3k6i5/flashtext
2 $ cd flashtext/docs
3 $ pip install sphinx
4 $ make html
5 $ # open _build/html/index.html in browser to view it locally

```

Why not Regex? It's a custom algorithm based on Aho-Corasick algorithm and Trie Dictionary.

Benchmark Time taken by FlashText to find terms in comparison to Regex.

<https://thepracticaldev.s3.amazonaws.com/i/xruf50n6z1r37ti8rd89.png> Time taken by FlashText to replace terms in comparison to Regex.

<https://thepracticaldev.s3.amazonaws.com/i/k44ghwp8o712dm58debj.png> Link to code for benchmarking the Find Feature and Replace Feature.

The idea for this library came from the following StackOverflow question.

Citation The original paper published on FlashText algorithm.

## 4 flashtext API

### 4.1 Import and initialize module

```

1 >>> from flashtext import KeywordProcessor
2 >>> keyword_processor = KeywordProcessor()
3 >>> # if match has to be case sensitive
4 >>> keyword_processor = KeywordProcessor(case_sensitive=True)

```

## 4.2 Add Keywords to module

```
1 >>> keyword_processor.add_keyword('Big Apple', 'New York')
2 >>> keyword_processor.add_keyword('Bay Area')
3 Extract keywords
4 >>> keywords_found = keyword_processor.extract_keywords('I love Big Apple and Bay
↳ Area.')
5 >>> keywords_found
6 >>> ['New York', 'Bay Area']
```

## 4.3 Replace keywords

```
1 >>> keyword_processor.add_keyword('New Delhi', 'NCR region')
2 >>> new_sentence = keyword_processor.replace_keywords('I love Big Apple and new
↳ delhi.')
3 >>> new_sentence
4 >>> 'I love New York and NCR region.'
```

## 4.4 Add keywords from File

```
1 >>> # Option 1: keywords.txt content
2 >>> # java_2e=>java
3 >>> # java programing=>java
4 >>> # product management=>product management
5 >>> # product management techniques=>product management
6 >>> # Option 2: keywords.txt content
7 >>> # java
8 >>> # python
9 >>> # c++
10 >>> keyword_processor.add_keyword_from_file('keywords.txt')
```

## 4.5 Add keywords from dict

```
1 >>> keyword_dict = {
2     "java": ["java_2e", "java programing"],
3     "product management": ["PM", "product manager"]
```

```
4     }  
5 >>> keyword_processor.add_keywords_from_dict(keyword_dict)
```

## 4.6 Add keywords from list

```
1 >>> keyword_processor.add_keywords_from_list(["java", "python"])
```