# 1、默认的分词器

standard

standard tokenizer：以单词边界进行切分
standard token filter：什么都不做。。。。
lowercase token filter：将所有字母转换成小写
stop token filter（默认被禁用）：移除停用词，比如：a the it 等等

# 2、修改分词器的设置

启用 english 停用词 token filter

```
1  PUT /my_index1                          // 启用停用词 token filter
2  {
3    "settings": {
4      "analysis": {
5        "analyzer": {
6          "es_std":{
7            "type":"standard",
8            "stopwords":"_english_"
9          }
10       }
11     }
12   }
13 }
14
15 GET /my_index1/_analyze
16 {
17   "analyzer": "es_std",
18   "text": ["i have a dog,and i love my dog"]
19 }
```

# 3、定制自己的分词器

```
1  PUT /my_index2
2  {
```

```
 3    "settings": {
 4      "analysis": {
 5        "char_filter": {
 6          "&_to_and":{
 7            "type":"mapping",
 8            "mappings":["&=> and "]
 9          }
10        },
11        "filter": {
12          "my_stopwords":{
13            "type":"stop",
14            "stopwords":["the", "a"]
15          }
16        },
17        "analyzer": {
18          "my_analyer":{
19            "type":"custom",
20            "char_filter":["html_strip","&_to_and"],
21            "tokenizer":"standard",
22            "filter":["lowercase","my_stopwords"]
23          }
24        }
25      }
26    }
27  }
28
29  GET /my_index2/_analyze
30  {
31    "analyzer": "my_analyer",
32    "text": ["i love my dad&mom,<a>pilipala</a>, the dog is so CUTE"]
33  }
```