In [2]:
```python
# Tokenise into sentences (returns a list):
from nltk.tokenize import sent_tokenize, word_tokenize

EXAMPLE_TEXT = "Hello Mr. Jones, how are you doing today? The weather is great. The sky is pinkish-blue. You shouldn't need an umbrella."
print(sent_tokenize(EXAMPLE_TEXT))
```

['Hello Mr. Jones, how are you doing today?', 'The weather is great.', 'The sky is pinkish-blue.', 'You shouldn't need an umbrella.']

In [3]:
```python
# Tokenise into words (returns a list):
from nltk.tokenize import sent_tokenize, word_tokenize

EXAMPLE_TEXT = "Hello Mr. Jones, how are you doing today? The weather is great. The sky is pinkish-blue. You shouldn't need an umbrella."
print(word_tokenize(EXAMPLE_TEXT))
```

['Hello', 'Mr.', 'Jones', ',', 'how', 'are', 'you', 'doing', 'today', '?', 'The', 'weather', 'is', 'great', '.', 'The', 'sky', 'is', 'pinkish-blue', '.', 'You', 'shouldn', ''', 't', 'need', 'an', 'umbrella', '.']

In [4]:
```python
# STOP WORDS - very common words that have little or no meaning
# thus should not be stored or processed
# do not filter out if PHRASE SEARCHING is required, though!
from nltk.corpus import stopwords
print(set(stopwords.words('english')))
```

{"that'll", 'an', 'mustn', 'about', 'after', 'been', 'on', 'yourself', "should've", 'doesn', 'that', 'have', 'no', 'my', 'hasn', "haven't", 'had', 'only', 'because', "couldn't", "weren't", 'against', "mightn't", 'can', 'was', "needn't", 'ain', 'being', 'those', 'from', 'having', 'doing', 'mightn', 'hers', 'haven', "doesn't", 'before', 'while', 'we', 'for', 'whom', 'yourselves', 'herself', 'some', 'than', 'to', 'shouldn', 'itself', 'now', 'what', 'which', 'off', 'do', "you're", 'its', 'and', "she's", 'out', 'over', 'when', 'few', 'a', 'don', 'she', 'just', 'couldn', 'him', 't', 'very', 'as', 'll', 'her', 'if', 'me', 'here', 'through', 'under', 'wouldn', 'myself', 'most', "shouldn't", 'above', 'isn', "didn't", 'in', 'own', 'not', 'how', 'of', 'you', "wouldn't", 'both', 'these', "you've", "hasn't", 'other', 'does', 'themselves', 'he', 'once', 'each', 'then', 'more', 'but', 'until', "it's", 've', "mustn't", 'yours', 'am', 'did', 'd', 'nor', 'their', 'too', 'has', 's', 'i', 'between', 'your', 'any', 'during', 'ours', "isn't", 'all', 'below', 'down', "you'd", 'himself', 'or', 'o', 'at', 'where', 're', 'with', 'should', 'so', 'there', 'who', 'theirs', 'same', "don't", 'it', 'wasn', "won't", 'the', 'weren', 'won', 'are', 'shan', "shan't", 'needn', 'y', "aren't", 'be', 'why', 'didn', 'into', 'up', 'ma', 'aren', 'by', 'his', 'our', "you'll", 'them', 'is', "wasn't", 'were', 'will', 'such', 'they', 'this', 'm', 'hadn', 'again', 'ourselves', "hadn't", 'further'}

In [5]:
```python
from nltk.corpus import stopwords
print(set(stopwords.words('german')))
```

{'an', 'damit', 'bin', 'haben', 'solchem', 'anderes', 'deinem', 'welcher', 'welche', 'also', 'gewesen', 'jenem', 'alle', 'soll', 'zu', 'einem', 'zum', 'und', 'wieder', 'habe', 'dann', 'dieselben', 'warst', 'solchen', 'der', 'allem', 'was', 'daß', 'während', 'von', 'dazu', 'das', 'hier', 'dem', 'eures', 'viel', 'einen', 'derselben', 'seine', 'sollte', 'aus', 'einiges', 'ihn', 'sein', 'derer', 'dir', 'eurer', 'meinen', 'seines', 'ist', 'eurem', 'als', 'andern', 'jenes', 'mancher', 'solcher', 'dort', 'jedem', 'anderer', 'im', 'meinem', 'derselbe', 'diesem', 'eine', 'andere', 'doch', 'indem', 'zwar', 'einig', 'dies', 'manches', 'kann', 'auf', 'ander', 'ich', 'keiner', 'eines', 'man', 'allen', 'da', 'sondern', 'denn', 'wir', 'unsere', 'ihrer', 'nichts', 'weil', 'euren', 'manchem', 'denselben', 'sehr', 'ihrem', 'einer', 'meines', 'in', 'jedes', 'musste', 'sonst', 'dieselbe', 'hinter', 'den', 'wo', 'jeder', 'ihnen', 'jeden', 'euer', 'unter', 'wollen', 'sich', 'einigem', 'welches', 'jede', 'ohne', 'einigen', 'welchem', 'können', 'gegen', 'etwas', 'jenen', 'manche', 'unseren', 'diesen', 'deiner', 'meiner', 'anderen', 'die', 'ins', 'machen', 'werde', 'jener', 'anderr', 'wie', 'bei', 'würde', 'mit', 'würden', 'einmal', 'am', 'mein', 'zur', 'ihm', 'seinen', 'dein', 'desselben', 'muss', 'wollte', 'deinen', 'bist', 'einiger', 'ihres', 'bis', 'deines', 'hatte', 'seiner', 'sie', 'sind', 'kein', 'werden', 'jetzt', 'könnte', 'dasselbe', 'einige', 'es', 'demselben', 'oder', 'diese', 'hin', 'mich', 'keinen', 'jene', 'vom', 'auch', 'nun', 'dieser', 'um', 'dieses', 'solches', 'keinem', 'so', 'uns', 'seinem', 'hatten', 'zwischen', 'aber', 'anderm', 'mir', 'ihre', 'solche', 'weiter', 'dich', 'keines', 'unseres', 'unserem', 'anderem', 'vor', 'nach', 'deine', 'hab', 'war', 'unser', 'noch', 'selbst', 'euch', 'ob', 'keine', 'welchen', 'meine', 'dessen', 'ihr', 'wenn', 'ihren', 'aller', 'manchen', 'über', 'alles', 'du', 'waren', 'er', 'für', 'hat', 'wird', 'eure', 'nur', 'des', 'will', 'weg', 'anders', 'wirst', 'nicht', 'ein', 'durch'}

In [2]:
```python
import nltk
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize

example_sent = ("This is a sample sentence, showing off the stop words filtration")

stop_words = set(stopwords.words('english'))
word_tokens = word_tokenize(example_sent)
filtered_sentence = [w for w in word_tokens if not w in stop_words]

filtered_sentence
print(filtered_sentence)
```

['This', 'sample', 'sentence', ',', 'showing', 'stop', 'words', 'filtration']

In [4]:

```python
# Alternative to list comprehension approach

from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize

example_sent = ("In diesem Satz, so wie bei uns, sind die Beispielwörter")

stop_words = set(stopwords.words('german'))
word_tokens = word_tokenize(example_sent)
filtered_sentence2 = []
for w in word_tokens:
    if w not in stop_words:
        filtered_sentence2.append(w)

filtered_sentence2
print(filtered_sentence2)
```

```
['In', 'Satz', ',', ',', 'Beispielwörter']
```