

# Dive into natural language processing with Python and NLTK

Andrii Soldatenko  
[@a\\_soldatenko](https://twitter.com/a_soldatenko)

# Agenda:

- Who am I?
- What is Natural language processing?
- Natural Language Toolkit (NLTK)
- Named-entity recognition (NER)
- Search speed
- What's next?

# Andrii Soldatenko

- Senior Python Developer at
- CTO in Ethoos.com
- Speaker at many PyCons and Python meetups
- blogger at <https://asoldatenko.com>



# Preface



# Siri



# Text Search

→ cpython time grep -r -i 'OrderedDict' .

grep -r -i 'OrderedDict' **2.35s user 0.10s system 97% cpu 2.510 total**

→ cpython time ack OrderedDict

ack OrderedDict **1.74s user 0.14s system 96% cpu 1.946 total**

→ cpython time pss OrderedDict

pss OrderedDict **0.85s user 0.09s system 96% cpu 0.983 total**

→ cpython time pt OrderedDict

pt OrderedDict **0.14s user 0.10s system 462% cpu 0.051 total**



Processor 2.5 GHz Intel Core i7  
Memory 16 GB 1600 MHz DDR3

# Natural Language Toolkit



# NLTK: Crash Course



# Natural Language Toolkit REST Api

```
→ curl -d "text=python" \
> http://text-processing.com/api/stem/
{
    "text": "python"
}
```

# Tokenization



# Tokenization: Example

```
>>> from nltk.tokenize import TreebankWordTokenizer  
  
>>> tokenizer = TreebankWordTokenizer()  
  
>>> tokenizer.tokenize('Software engineering  
conference')  
['Software', 'engineering', 'conference']
```

# Tokenization: Oops :(

```
>>> from nltk.tokenize import  
TreebankWordTokenizer  
>>> tokenizer = TreebankWordTokenizer()  
>>> tokenizer.tokenize("Can't tokenize  
this")  
[ 'Ca', "n't", 'tokenize', 'this' ]
```

# Tokenization: Oops :( 2

```
>>> from nltk.tokenize import  
WordPunctTokenizer  
>>> tokenizer = WordPunctTokenizer()  
>>> tokenizer.tokenize("Can't tokenize  
this")  
[ 'Can', "'", "'", 't', 'tokenize', 'this' ]
```

# Tokenization: Regular Expression

```
>>> from nltk.tokenize import  
RegexpTokenizer  
>>> tokenizer = RegexpTokenizer("[\w']+")  
>>> tokenizer.tokenize("Can't tokenize  
this")  
[ "Can't", 'tokenize', 'this' ]
```

When your regular expression  
returns what you expect:)



**WordNet - a lexical  
database for English  
language**

# Download Data

```
→ python -v
```

```
Python 3.5.2
```

```
→ python -m nltk.downloader all
```

```
[nltk_data] Downloading collection 'all'
```

```
[nltk_data]
```

```
[nltk_data] | Downloading package abc to /home/vagrant/nltk_data...  
[nltk_data] | Unzipping corpora/abc.zip.
```

```
[nltk_data]
```

```
[nltk_data] | Downloading package alpino to
```

```
...
```

```
→ python -m nltk.downloader punkt
```

```
[nltk_data] Downloading package punkt to...
```

```
[nltk_data]
```

```
[nltk_data] | Unzipping tokenizers/punkt.zip
```

# Stop Words

until      not in  
few or him ours it off  
does their your my he be  
as do we i ourselves a  
she myself am  
through what himself them  
and

# StopWords: total count

```
→ cat /home/vagrant/nltk_data/corpora/  
stopwords/english | wc -l  
153
```

```
→ cat /usr/share/postgresql/9.5/tsearch_data/  
english.stop | wc -l  
127
```

# StopWords: lang count

```
→ ls /home/vagrant/nltk_data/corpora/  
stopwords/ | grep -v README | wc -l
```

15

# StopWords: diff

```
→ diff /home/vagrant/nltk_data/corpora/  
stopwords/english /usr/share/postgresql/9.5/  
tsearch_data/english.stop
```

```
128,153d127
```

```
< d  
< ll  
< m  
< o  
< re  
< ve  
< y  
< ain  
< aren
```

# StopWords: Example

```
>>> from nltk.corpus import stopwords  
>>> english_stops =  
set(stopwords.words('english'))  
>>> words = ["Can't", "tokenize", "this"]  
>>> [word for word in words if word not in  
english_stops]  
["Can't tokenize"]
```



# Improve Stop words

```
import re

replacement_patterns = [
    (r>won't', 'will not'),
    (r>can't', 'can not'),
    ...
]

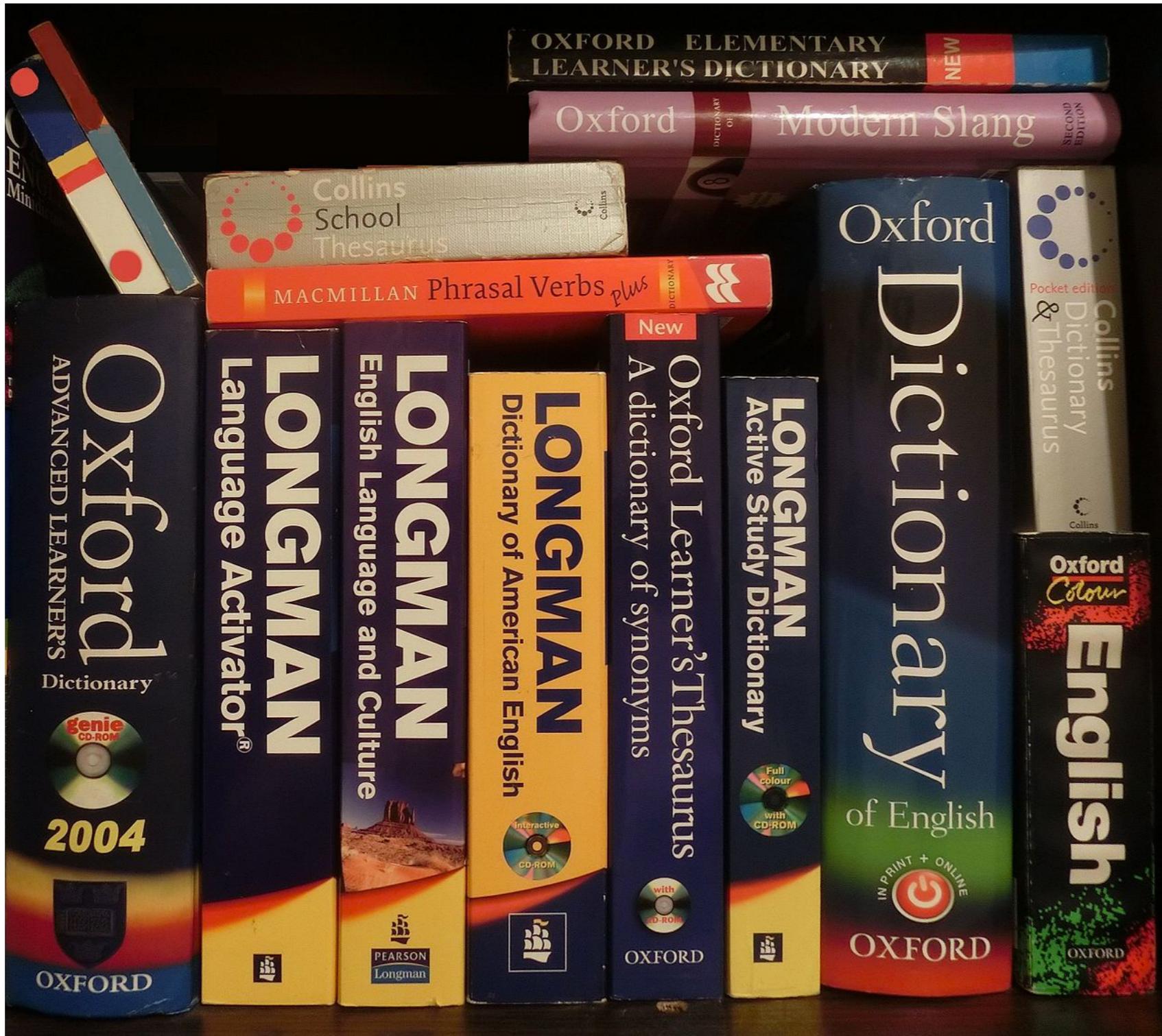
class RegexpReplacer(object):
    def __init__(self, patterns=replacement_patterns):
        self.patterns = [(re.compile(regex), repl) for (regex,
repl) in
                           patterns]

    def replace(self, text):
        s = text
        for (pattern, repl) in self.patterns:
            s = re.sub(pattern, repl, s)
        return s
```

# Improve Stop words

```
>>> english_stops = set(stopwords.words('english'))
>>> sentence = "Can't tokenize this"
>>> replacer = RegexpReplacer() ← Custom class
>>> tokenizer = WordPunctTokenizer()
>>>
>>> words =
tokenizer.tokenize(replacer.replace(sentence.lower()))
>>> [word for word in words if replacer.replace(word)
not in english_stops]
['tokenize']
```

# Collocations



# Collocations:Example

```
>>> words = [w.lower() for w in  
webtext.words('grail.txt')]  
>>> bcf = BigramCollocationFinder.from_words(words)  
>>> results =  
bcf.nbest(BigramAssocMeasures.likelihood_ratio, 10)  
>>> pprint.pprint([' '.join(r) for r in results])  
['black knight',  
'clop clop',  
'head knight',  
'mumble mumble',  
'squeak squeak',  
'saw saw',  
'holy grail',  
'run away',  
'french guard',  
'cartoon character']
```

Full example: <http://bit.ly/2c1Rzkv>

# Stemming



<https://tartarus.org/martin/PorterStemmer/def.txt>

# Porter stemming algorithm

```
>>> from nltk.stem import PorterStemmer
>>> stemmer = PorterStemmer()
>>> stemmer.stem('cooking')
'cook'
>>> stemmer.stem('cookery')
'cookeri'
```

source code of algorithm: nltk.stem.porter.PorterStemmer

# Lancaster stemming algorithm

```
>>> from nltk.stem import  
LancasterStemmer  
>>> stemmer = LancasterStemmer()  
>>> stemmer.stem('cooking')  
'cook'  
>>> stemmer.stem('cookery')  
'cookery'
```

source code: nltk.stem.lancaster.LancasterStemmer

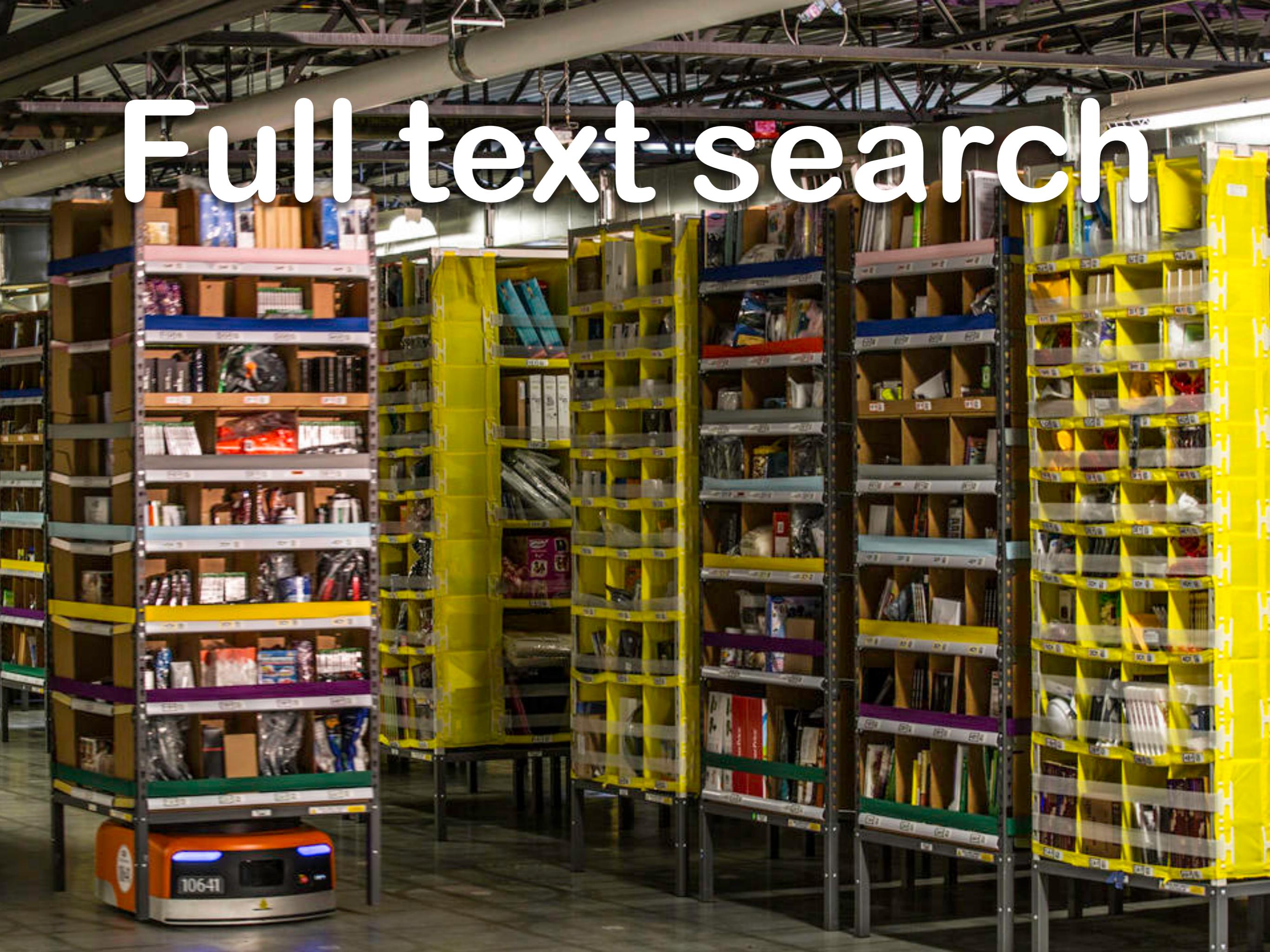
# Lemmatizing



# Lemmatizing

```
>>> from nltk.stem import (
...     PorterStemmer,
...     WordNetLemmatizer
... )
>>>
>>> stemmer = PorterStemmer()
>>> stemmer.stem('believes')
'believ'
>>>
>>> lemmatizer = WordNetLemmatizer()
>>> lemmatizer.lemmatize('believes')
'belief'
```

# Full text search



# Search index

## Symbols

32gb Heap boundary, 642

## A

ACID transactions, 545, 556

action, in bulk requests, 57, 69

ad hoc searches, 15

aggregations, 20, 417

  aggs parameter, 424

  and analysis, 483

  approximate, 457

  cardinality, 458

  percentiles, 462

basic example

  adding a metric, 426

  adding extra metrics, 429

  buckets nested in other buckets, 427

buckets, 419

  combining buckets and metrics, 420

  metrics, 420

limiting memory usage, 487

  fielddata circuit breaker, 490

  fielddata size, 488

  moitoring fielddata, 489

managing efficient memory usage, 507

nested, 567

  reverse\_nested aggregation, 568

operating alongside search requests, 418

preventing combinatorial explosions, 500

  depth-first versus breadth-first, 502

returning empty buckets, 439

scoping, 445

  global bucket, 447

Significant Terms, 471

significant terms

# Simple sentences

1. The quick brown fox jumped over the lazy dog
2. Quick brown foxes leap over lazy dogs in summer

# Inverted index

1	Term	Doc_1	Doc_2
2		-----	
3	Quick		X
4	The	X	
5	brown	X	X
6	dog	X	
7	dogs		X
8	fox	X	
9	foxes		X
10	in		X
11	jumped	X	
12	lazy	X	X
13	leap		X
14	over	X	X
15	quick	X	
16	summer		X
17	the	X	
18		-----	

# Inverted index

1	Term	Doc_1	Doc_2
2	-----		
3	brown		X   X
4	quick		X
5	-----		
6	Total	2	1

# Inverted index: normalization

Term	Doc_1	Doc_2
<hr/>		
Quick	X	X
The	X	
brown	X	X
dog	X	
dogs	X	
fox	X	
foxes	X	
<b>in</b>	X	
jumped	X	
lazy	X   X	
leap	X	
over	X   X	
quick	X	
summer	X	
the	X	
<hr/>		

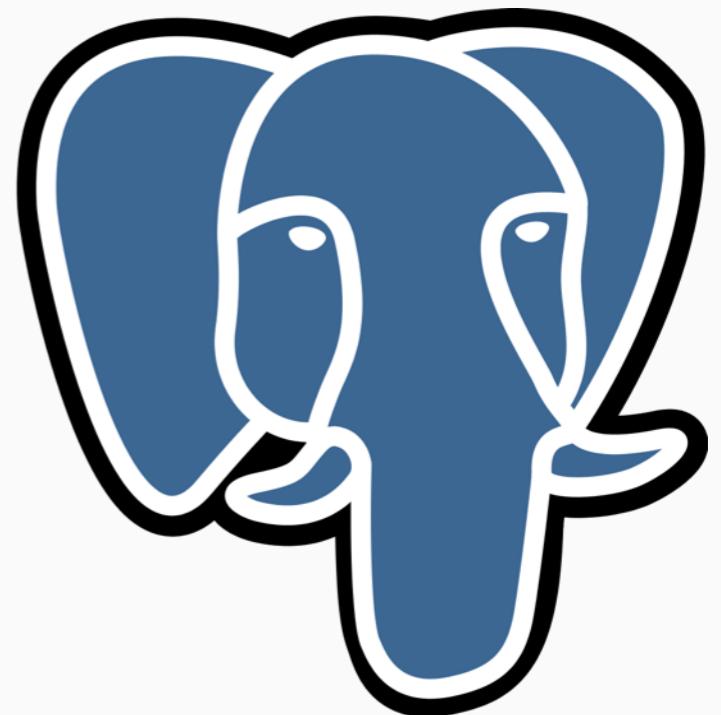
Diagram illustrating the normalization of an inverted index. A green arrow points from the original index on the left to a normalized version on the right. The normalized index removes terms that appear in only one document and renames terms that appear in both documents.

Term	Doc_1	Doc_2
brown	X   X	
dog	X   X	
fox	X   X	
<b>in</b>	X	
jump	X   X	
lazy	X   X	
over	X   X	
quick	X   X	
summer	X	
the	X   X	

# What is the best full text search engine in Python?



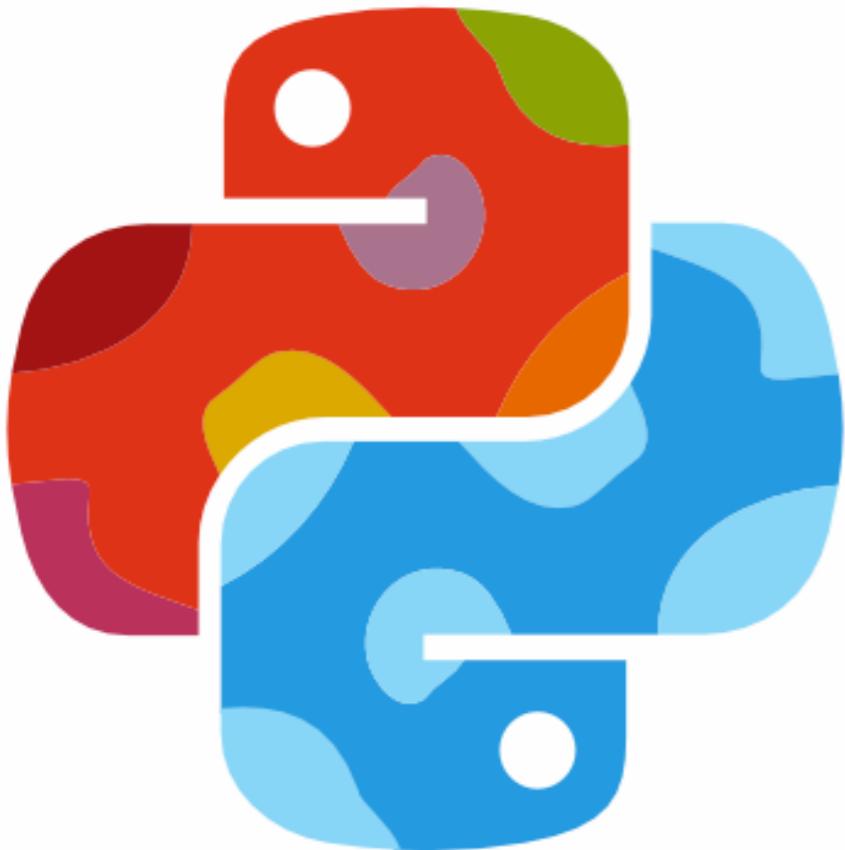
elastic



<https://asoldatenko.com/europython-2016.html>

# NER

<https://tartarus.org/martin/PorterStemmer/def.txt>



# EUROPYTHON

## 2016

Bilbao, 17-24 July

My talk about search engines you can find:

<https://asoldatenko.com/europython-2016.html>

# 1 million music Artists

Evie Tamala

Jean-Pierre Martin

Deejay One

wecamewithbrokenteeth

The Blackbelt Band

Giant Tomo

Decoding Jesus

Elvin Jones & Jimmy Garrison Sextet

Infester

...

David Silverman

Aili Teigmo

## Performance

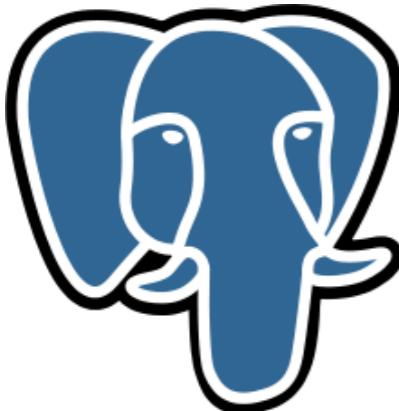
## Database size



elastic

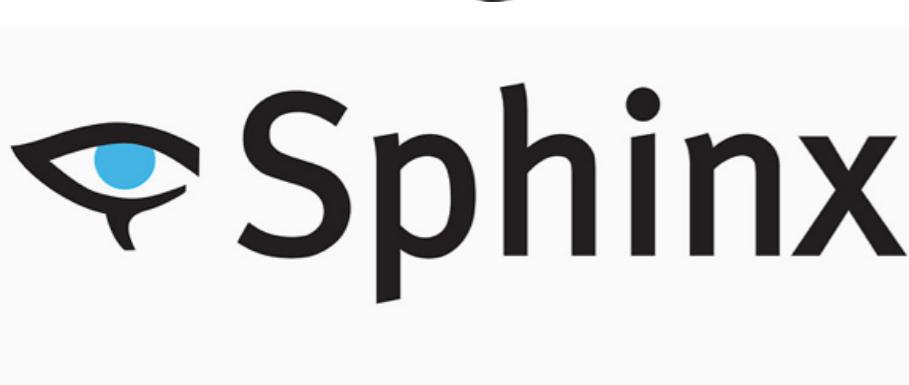
9 ms

~ 1 million records



4 ms

~ 1 million records



6 ms

~ 1 million records



~2 s

~ 1 million records

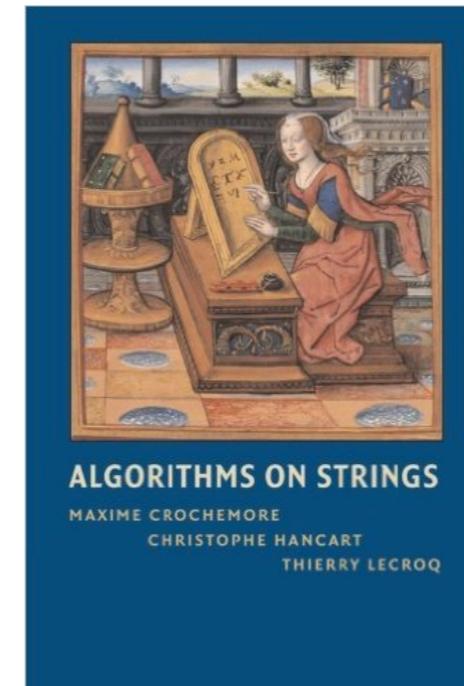
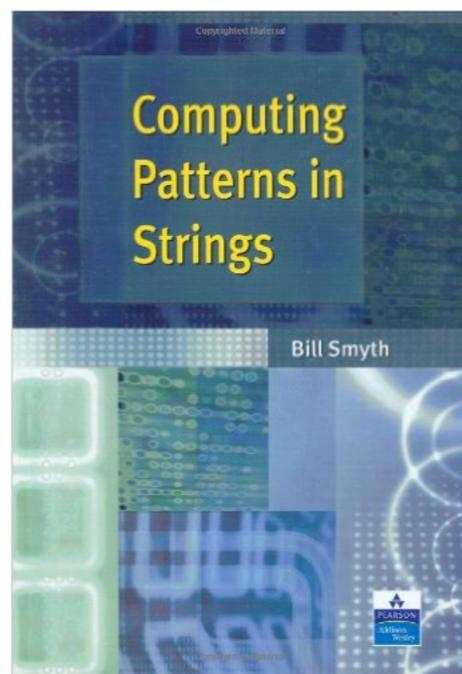
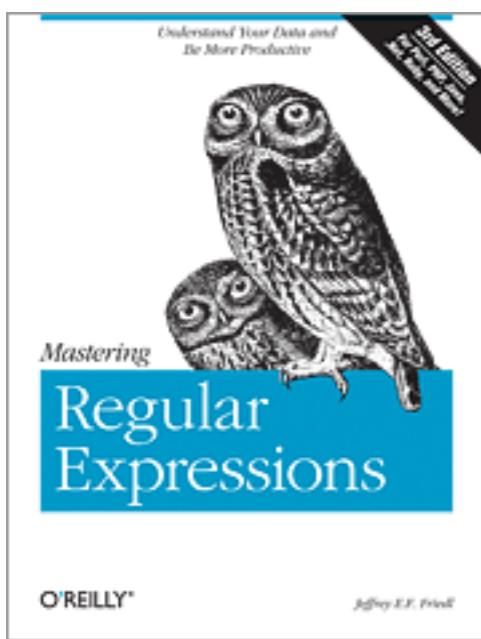
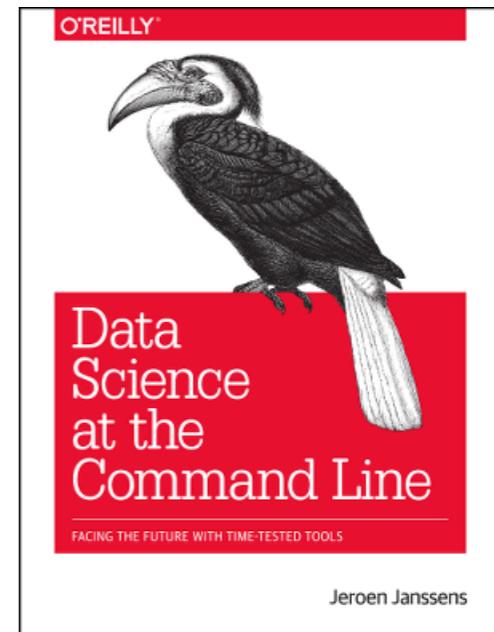
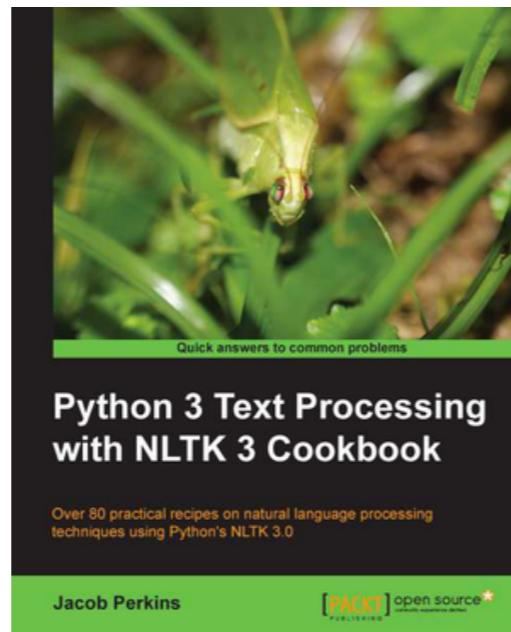
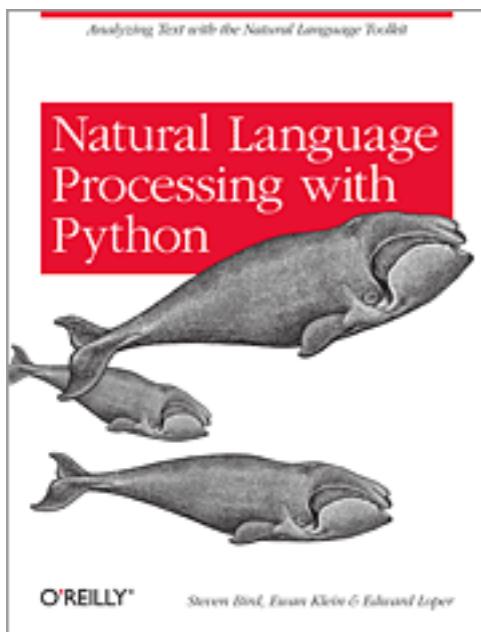
# Named-entity recognition

- Jim bought 300 shares of Acme Corp. in 2006.
- [Jim]Person bought 300 shares of [Acme Corp.]Organization in [2006]Time.

# Books



# Books



# References:

[http://www.nltk.org/book\\_1ed/](http://www.nltk.org/book_1ed/)

<http://bit.ly/2bKXnlz>

# Slides

<https://asoldatenko.com/pdfs/se2016.pdf>

<http://bit.ly/2bJAuJb>

# Thank You



@a\_soldatenko



andrii.soldatenko@gmail.com



Hire the top 3% of freelance developers

<http://bit.ly/21lxQ01>



/[Q|q]uestions/

