

---

# **spaCy Documentation**

***Release 0.1***

**Matthew Honnibal**

August 20, 2014



---

Contents

---

<b>1</b>	<b>Python API</b>	<b>1</b>
1.1	Cheat Sheet . . . . .	1
<b>2</b>	<b>Cython API</b>	<b>3</b>
2.1	Cheat Sheet . . . . .	3
<b>3</b>	<b>Adding a Language</b>	<b>5</b>
<b>4</b>	<b>Overview</b>	<b>7</b>
<b>5</b>	<b>Pros and Cons</b>	<b>9</b>
<b>6</b>	<b>Installation</b>	<b>11</b>



---

**Python API**

---

## 1.1 Cheat Sheet



---

**Cython API**

---

## 2.1 Cheat Sheet



---

## **Adding a Language**

---



---

## Overview

---

spaCy is a tokenizer for natural languages, tightly coupled to a global vocabulary store.

Instead of a list of strings, spaCy returns references to lexical types. All of the string-based features you might need are pre-computed for you:

```
>>> from spacy import en
>>> example = u"Apples aren't oranges..."
>>> apples, are, nt, oranges, ellipses = en.tokenize(example)
>>> en.is_punct(ellipses)
True
>>> en.get_string(en.word_shape(apples))
'XXXX'
```

You also get lots of distributional features, calculated from a large sample of text:

```
>>> en.prob_of(are) > en.prob_of(oranges)
True
>>> en.can_noun(are)
False
>>> en.is_oft_title(apples)
False
```



### Pros and Cons

---

Pros:

- All tokens come with indices into the original string
- Full unicode support
- Extensible to other languages
- Batch operations computed efficiently in Cython
- Cython API
- numpy interoperability

Cons:

- It's new (released September 2014)
- Higher memory usage (up to 1gb)
- More conceptually complicated
- Tokenization rules expressed in code, not as data



---

## Installation

---

Installation via pip:

```
pip install spacy
```

From source, using virtualenv:

```
$ git clone http://github.com/honnibal/spaCy.git
$ cd spaCy
$ virtualenv .env
$ source .env/bin/activate
$ pip install -r requirements.txt
$ fab make
$ fab test
```