

◆ Member-only story

# In-Depth spaCy Tutorial For Beginners in NLP

Learn the Scikit-learn of Natural Language Processing



Bex T. · [Follow](#)



Published in Towards Data Science · 9 min read · Jun 9, 2022



128



Photo by [Anni Roenkae](#)

## Introduction

No, we won't be building language models with billions of parameters today. We will start smaller and learn about the basics of NLP with spaCy. We will closely look at how the library works and how you can use it to solve beginner/intermediate NLP problems with ease.

The post is already long, so I'll cut it here and jump to the meat of the article.

[Open in app ↗](#)



Search



Write



~~standard with vast features to solve many NLP tasks with state-of-the-art speed, accuracy, and performance.~~

At its core are pipelines, which you can think of as language-specific models already trained on millions of text instances.

It is also at the head of the [spaCy ecosystem](#) that includes dozens of libraries and tools such as Prodigy, Forte, displaCy, explacy, ADAM, Coreferee, etc.

spaCy can also shake hands with custom models from TensorFlow, PyTorch, and other frameworks.

Right now, spaCy supports 66 languages as separate pipelines, and new languages are being added slowly.

## Basics of spaCy

Before seeing how spaCy works, let's install it:

```
pip install -U spacy
```

spaCy has three pipelines for English, with varying sizes and functionality for complex tasks. In this tutorial, we will only need to install the small and medium pipelines, but I included the large one as well for completeness:

```
1 python -m spacy download en_core_web_sm # 12 MB
2 python -m spacy download en_core_web_md # 31 MB
3 python -m spacy download en_core_web_lg # 382 MB
```

11801.sh hosted with ❤️ by GitHub

[view raw](#)

After importing spaCy, we need to load one of the pipelines we just installed. For now, we will load the small one and store it to `nlp`:

```
1 import spacy
2
3 txt = "The tallest living man is 37-year-old Sultan Kosen, from Turkey,"
4     " who is 8 feet, 2.8 inches, who set the record in 2009."
5
6 # Create the Language object
7 nlp = spacy.load("en_core_web_sm")
8 nlp
9 <spacy.lang.en.English at 0x1f33e3af550>
```

11802.py hosted with ❤️ by GitHub

[view raw](#)

It is a convention to name any loaded language models `nlp` in the spaCy ecosystem. This object can now be called on any text to start information extraction:

```
# Create the Doc object
doc = nlp(txt)
```

The `doc` object is also a convention, and now it is already filled with extra information about the given text's sentences and words.

In general, the `doc` object is just an iterator:

```
1 for token in doc[:5]:
2     print(token)
3
4 [OUT]:
5
6 The
7 tallest
8 living
9 man
10 is
```

11803.py hosted with ❤ by GitHub

[view raw](#)

You can use slicing or indexing notations to extract individual tokens:

```
>>> type(token)
spacy.tokens.token.Token
>>> len(doc)
```

31

*Tokenization is splitting sentences into words and punctuation. A single token can be a word, a punctuation or a noun chunk, etc.*

If you extract more than one token, then you have a span object:

```
>>> span = doc[:5]
>>> type(span)
spacy.tokens.span.Span
>>> span.text
'The tallest living man is'
```

spaCy is also built for memory efficiency. That's why both `token` and `span` objects are just views of the `doc` object. There is no duplication.

The pre-trained English pipeline and many other pipelines have language-specific rules for tokenization and extracting their lexical attributes. Here are 6 of such attributes:

```

1  print("Index:    ", [token.i for token in doc[3:10]])
2  print("Text:     ", [token.text for token in doc[3:10]])
3  print("is_alpha: ", [token.is_alpha for token in doc[3:10]])
4  print("is_punct: ", [token.is_punct for token in doc[3:10]])
5  print("like_num: ", [token.like_num for token in doc[3:10]])
6  print("Base word:", [token.lemma_ for token in doc[3:10]])
7
8  [OUT]:
9
10 Index:    [3, 4, 5, 6, 7, 8, 9]
11 Text:     ['man', 'is', '37', '-', 'year', '-', 'old']
12 is_alpha: [True, True, False, False, True, False, True]
13 is_punct: [False, False, False, True, False, True, False]
14 like_num: [False, False, True, False, False, False, False]
15 Base word: ['man', 'be', '37', '-', 'year', '-', 'old']

```

11804.py hosted with ❤️ by GitHub

[view raw](#)

Some interesting attributes are `lemma_`, which returns the base word stripped from any suffixes, prefixes, tense, or any other grammatical attributes, and the `like_num` which recognizes both literal and lettered numbers.

You will be spending most of your time on these four objects — `nlp`, `doc`, `token` and `span`. Let's take a closer look at how they are related.

## Architecture and core data structures

Let's start again with the `nlp`, which is a Language object under the hood:

```
1 import spacy
2
3 nlp = spacy.load("en_core_web_md")
4
5 >>> type(nlp)
6 spacy.lang.en.English
```

11805.py hosted with ❤️ by GitHub

[view raw](#)

Language objects are pre-trained on millions of text instances and labels and loaded into spaCy with their binary weights. These weights allow you to perform various tasks on new datasets without worrying about the hairy details.

As I mentioned earlier, spaCy has fully-trained pipelines for 22 languages, some of which you can see below:

```
1 nlp = spacy.load("es_core_news_sm") # Spanish
2 nlp = spacy.load("ru_core_news_sm") # Russian
3 nlp = spacy.load("zh_core_web_sm") # Chinese
4 nlp = spacy.load("de_core_news_sm") # German
```

11806.py hosted with ❤️ by GitHub

[view raw](#)

For other +40 languages, spaCy only offers basic tokenization rules, and other functionality is being slowly integrated with community effort.

It is also possible to load the language models directly from the `lang` submodule:

```

1  from spacy.lang.en import English
2  from spacy.lang.es import Spanish
3
4  nlp = English()
5
6  >>> type(nlp)
7  spacy.lang.en.English
8
9  txt = """The original name for the search engine Google was Backrub.
10         It was renamed Google after the googol,
11         which is the number one followed by 100 zeros."""
12
13 doc = nlp(txt)

```

11807.py hosted with ❤ by GitHub

[view raw](#)

After processing a text, words and punctuation are stored in the vocabulary object of `nlp`:

```

>>> type(nlp.vocab)
spacy.vocab.Vocab

```

This `Vocab` is shared between documents, meaning it stores all new words from all docs. In contrast, the `doc` object's vocabulary only contains the words from the `txt`:

```

>>> type(doc.vocab)
spacy.vocab.Vocab

```

Internally, spaCy communicates in hashes to save memory and has a two-way lookup table called `StringStore`. You can get the hash of a string or get the string if you have the hash:

```
>>> type(nlp.vocab.strings)
spacy.strings.StringStore
>>> nlp.vocab.strings["google"]
1988622737398120358
>>> nlp.vocab.strings[1988622737398120358]
'google'
```

When tokens go into the `Vocab`, they lose all their context-specific information. So, when you look up words from the vocab, you are looking up lexeme S:

```
>>> lexeme = nlp.vocab["google"]
>>> type(lexeme)
spacy.lexeme.Lexeme
```

Lexemes don't contain context-specific information like part-of-speech tags, morphological dependencies, etc. But they still offer many lexical attributes of the word:

```
>>> print(lexeme.text, lexeme.orth, lexeme.is_digit)
google 1988622737398120358 False
```

*| orth attribute is for the hash of the lexeme.*

So, if you are looking at a word through the `doc` object, it is a token. If it is from a `Vocab`, it is a lexeme.

Now, let's talk more about the `doc` object.

```
1  txt = """Mosquitoes are the deadliest animal in the world:  
2      They kill more people than any other creature,  
3      due to the diseases they carry."""  
4  
5  doc = nlp(txt)  
6  
7  >>> type(doc)  
8  spacy.tokens.Doc
```

11808.py hosted with ❤️ by GitHub

[view raw](#)

Calling `nlp` object on text generates the `doc` along with its special attributes.

```
1  >>> doc.text  
2  'Mosquitoes are the deadliest animal in the world: \n'  
3  '      They kill more people than any other creature, '  
4  '\n      due to the diseases they carry.'  
5  
6  >>> len(doc)  
7  27
```

11809.py hosted with ❤️ by GitHub

[view raw](#)

You can create docs manually by importing the `Doc` class from `tokens` module:

```
1 from spacy.lang.en import English
2 from spacy.tokens import Doc
3
4 nlp = English()
5
6 words = ["I", "love", "Barcelona", "!"]
7 spaces = [True, True, False, False]
8
9 # Create the doc object manually
10 doc = Doc(nlp.vocab, words=words, spaces=spaces)
```

11810.py hosted with ❤ by GitHub

[view raw](#)

Doc requires three arguments - the vocabulary from `nlp`, a list of words and another list specifying if the words are followed by a space (including the last one). All tokens in `doc` have this information.

```
>>> len(doc)
4
>>> doc.text
'I love Barcelona!'
```

Spans are also a class of their own and expose a range of attributes, even though they are just a view of the `doc` object:

```

1  txt = """The hardest working muscle in your body is your heart:  

2      It pumps more than 2,000 gallons of blood a day  

3      and beats more than 2.5 billion times in a 70-year life span."""  

4  

5  doc = nlp(txt)  

6  

7  span = doc[:10]  

8  

9  >>> type(span)  

10 spacy.tokens.span.Span  

11  

12 >>> print(span.text)  

13 The hardest working muscle in your body is your heart  

14  

15 >>> print(span.start, span.end)  

16 0 10

```

11811.py hosted with ❤️ by GitHub

[view raw](#)

To create the span objects manually, pass the doc object and start/end indices of the tokens to the `Span` class:

```

1  from spacy.tokens import Span  

2  

3  span = Span(doc, 0, 10)  

4  

5  >>> span.text  

6  'The hardest working muscle in your body is your heart'

```

11812.py hosted with ❤️ by GitHub

[view raw](#)

## Named Entity Recognition (NER)

One of the most common tasks in NLP is predicting named entities, like people, locations, countries, brands, etc.

Performing NER is ridiculously easy in spaCy. After processing a text, just extract the `ents` attribute of the `doc` object:

```

1  txt = """Cleopatra wasn't actually Egyptian!
2      As far as historians can tell, Egypt's
3          famous femme fatal was actually Greek!.
4      She was a descendant of Alexander the Great's
5          Macedonian general Ptolemy"""
6
7  nlp = spacy.load("en_core_web_md")
8
9  doc = nlp(txt)
10
11 for ent in doc.ents:
12     print(f"{ent.text}: {ent.label_}")
13
14 [OUT]:
15
16 Cleopatra    PERSON
17 Egyptian     NORP
18 Egypt        GPE
19 Greek         NORP
20 Macedonian   NORP
21 Ptolemy      PERSON

```

11813.py hosted with ❤ by GitHub

[view raw](#)

Cleopatra is recognized as a PERSON, while Egypt is a Geo-political entity (GPE). To know the meaning of other labels, you can use the `explain` function:

```

1  >>> spacy.explain("GPE")
2  'Countries, cities, states'
3
4  >>> spacy.explain("NORP")
5  'Nationalities or religious or political groups'

```

11814.py hosted with ❤ by GitHub

[view raw](#)

Instead of printing text, you can use spaCy's visual entity tagger, available via displacy:

```
[30]: from spacy import displacy  
  
displacy.render(doc, style="ent")
```

```
Cleopatra PERSON wasn't actually Egyptian NORP !  
As far as historians can tell, Egypt GPE 's  
famous femme fatal was actually Greek NORP !.  
She was a descendant of Alexander the Great's  
Macedonian NORP general Ptolemy PERSON
```

Image by author

The image shows that Alexander the Great isn't recognized as a PERSON because it is not a common name. But no matter, we can label Alexander as a PERSON manually.

First, extract the full name as a span by giving it a label (PERSON):

```
1 from spacy.tokens import Span  
2  
3 alexander = Span(doc, 31, 34, label="PERSON")  
4 alexander  
5 Alexander the Great
```

11815.py hosted with ❤️ by GitHub

[view raw](#)

Then, update the `ents` list with the span:

```
doc.ents = list(doc.ents) + [alexander]
```

Now, displacy tags it as well:

The screenshot shows a Jupyter Notebook cell with the code `[62]: displacy.render(doc, style="ent")`. The output is a visual representation of the document's entities. The text is in white on a black background, and entities are highlighted with colored boxes. The entities and their spans are:

- Cleopatra PERSON
- wasn't actually Egyptian NORP !
- As far as historians can tell, Egypt GPE 's
- famous femme fatal was actually Greek NORP !.
- She was a descendant of Alexander the Great PERSON 's
- Macedonian NORP general Ptolemy PERSON

Image by author

You could've set the new entity with `set_ents` function as well:

```
# Leaves the rest of ents untouched  
doc.set_ents([alexander], default="unmodified")
```

## Predicting part-of-speech (POS) tags and syntactic dependencies

spaCy also offers a rich selection of tools for grammar analysis of a document. Lexical and grammatical attributes of tokens are given as attributes.

For example, let's take a look at each token's part-of-speech tag and its syntactic dependency:

```
1  txt = "The first footprints on the moon will remain there for a million years"
2
3  doc = nlp(txt)
4
5  print(
6      f"{'Text':<20} {'Part-of-speech':<20} "
7      f"{'Dependency':<20} {'Dependency text':<20}\n"
8  )
9  for token in doc:
10     print(f"{token.text:<20} {token.pos_:<20} "
11           f"{token.dep_:<20} {token.head.text:<20}")
12
13 [OUT]:
14
15 Text          Part-of-speech      Dependency      Dependency text
16
17 The            DET              det             footprints
18 first          ADJ              amod            footprints
19 footprints    NOUN             nsubj            remain
20 on              ADP              prep             footprints
21 the            DET              det              moon
22 moon           NOUN             pobj             on
23 will           AUX              aux              remain
24 remain         VERB             ROOT             remain
25 there          ADV              advmod            remain
26 for             ADP              prep             remain
27 a               DET              quantmod        million
28 million        NUM              nummod            years
29 years          NOUN             pobj             for
```

11816.py hosted with ❤ by GitHub

[view raw](#)

The output contains some confusing labels, but we can infer some of them like verbs, adverbs, adjectives, etc. Let's see the explanation for a few others:

```
1 pos_tags = ["DET", "AUX", "ADP"]
2 dep_tags = ["amod", "nsubj", "nummod"]
3
4 for pos in pos_tags:
5     print(pos, "-->", spacy.explain(pos))
6
7 for dep in dep_tags:
8     print(dep, "-->", spacy.explain(dep))
9
10 [OUT]:
11
12 DET --> determiner
13 AUX --> auxiliary
14 ADP --> adposition
15 amod --> adjectival modifier
16 nsubj --> nominal subject
17 nummod --> numeric modifier
```

11817.py hosted with ❤ by GitHub

[view raw](#)

The last column in the previous table represents word relations like "the first", "first footprints", "remain there", etc.

spaCy contains many more powerful features for linguistic analysis. As the last example, here is how you extract noun chunks:

```
1  txt = """The teddy bear is named after President Theodore Roosevelt.  
2      After he refused to shoot a captured black bear on a hunt,  
3      a stuffed-animal maker decided to create  
4      a bear and name it after the president."""  
5  
6  doc = nlp(txt)  
7  
8  for chunk in doc.noun_chunks:  
9      print(chunk.text)  
10  
11 [OUT]:  
12  
13 The teddy bear  
14 President Theodore Roosevelt  
15 he  
16 a captured black bear  
17 a hunt  
18 a stuffed-animal maker  
19 a bear  
20 it  
21 the president
```

1818.py hosted with ❤ by GitHub

[view raw](#)

Learn more about linguistic features from [this page](#) of the spaCy User Guide.

## Custom rule-based tokenization

Until now, spaCy had complete control over tokenization rules. But a language can have many culture-specific idiosyncrasies and edge-cases that don't fit spaCy's rules. For example, in an earlier example, we saw that "Alexander the Great" was missed as an entity.

```

1  txt = """Cleopatra wasn't actually Egyptian!
2      As far as historians can tell, Egypt's
3          famous femme fatal was actually Greek!.
4      She was a descendant of Alexander the Great's
5          Macedonian general Ptolemy"""

```

11819.py hosted with ❤️ by GitHub

[view raw](#)

If we process the same text again, spaCy regards the entity as three tokens. We need to tell it that titled names like Alexander the Great or Bran the Broken should be considered a single token rather than three because splitting them makes no sense.

Let's see how to do that by creating custom tokenization rules.

We will start by creating a pattern as a dictionary:

```

1  # Create a pattern
2  pattern = [
3      {"IS_ALPHA": True, "IS_TITLE": True},
4      {"IS_STOP": True},
5      {"IS_ALPHA": True, "IS_TITLE": True},
6  ]

```

11820.py hosted with ❤️ by GitHub

[view raw](#)

In spaCy, there is a large set of keywords you can use in combination to parse virtually any type of token pattern. For example, the above pattern is for a three-token pattern, with the first and last tokens being an alphanumeric text and the middle one being a stop word (like *the*, *and*, *or*, *etc.*). In other words, we are matching "Alexander the Great" without explicitly telling it to spaCy.

Now, we will create a `Matcher` object with this pattern:

```

1 from spacy.matcher import Matcher
2
3 # Init the matcher with the shared vocab
4 matcher = Matcher(nlp.vocab)
5
6 # Add the pattern to the matcher
7 matcher.add("TITLED_PERSON", [pattern])

```

11822.py hosted with ❤️ by GitHub

[view raw](#)

After processing the text, we call this matcher object on the `doc` object, which returns a list of matches. Each match is a tuple with three elements - match ID, start, and end:

```

1 # Process the text
2 doc = nlp(txt)
3
4 # Find all matches
5 matches = matcher(doc)
6
7 # Iterate over matches
8 for match_id, start, end in matches:
9     # Get the span
10    span = doc[start:end]
11    print(span.text)
12
13
14 [OUT]: Alexander the Great

```

11823.py hosted with ❤️ by GitHub

[view raw](#)

You can tweak the pattern in any you want. For example, you can use quantifiers like `OP` with reGex keywords or use reGex itself:

```

1 pattern = [
2     {"TEXT": {"REGEX": "[a-zA-Z]+"}},
3     {"IS_DIGIT": True, "OP": "?"}, # Match one or more times
4     {"DEP": "quantmod"}, # Match based on dependency
5     # etc.
6 ]

```

11824.py hosted with ❤️ by GitHub

[view raw](#)

You can learn more about custom rule-based matching from [here](#).

## Word vectors and semantic similarity

Another everyday use case of NLP is predicting semantic similarity. Similarity scores can be used in recommender systems, plagiarism, duplicate content, etc.

spaCy calculates semantic similarity using word vectors (explained below), which are available in the medium-sized model:

```

1 nlp = spacy.load("en_core_web_md")
2
3 doc1 = nlp("What a lukewarm sentiment.")
4 doc2 = nlp("What a short sentence.")
5
6 >>> doc1.similarity(doc2)
7 0.9200780919749721

```

11825.py hosted with ❤️ by GitHub

[view raw](#)

All `doc`, `token` and `span` objects have this `similarity` method:

```

1 token1 = doc1[-1]
2 token2 = doc2[2]
3
4 >>> token1.similarity(token2)
5 0.5222234129905701

```

11826.py hosted with ❤️ by GitHub

[view raw](#)

The three classes can be compared to each other as well, like a token to a span:

```

>>> doc1[0:2].similarity(doc[3])
0.8700238466262817

```

The similarity is calculated using word vectors, which are multi-dimensional mathematical representations of words. For example, here is the vector of the first token in the document:

```

1 array = doc1[0].vector
2
3 >>> array.shape
4 (300,)
5
6 >>> array[:10]
7 array([-0.73351,  0.41392, -0.4425, -0.29127, -0.096179,  0.097562,
8       0.13151, -0.51825,  0.10671,  2.4144], dtype=float32)

```

11827.py hosted with ❤️ by GitHub

[view raw](#)

## All about pipelines

Under the hood, language models aren't one pipeline but a collection:

```

1 nlp = spacy.load("en_core_web_sm")
2
3 >>> nlp.pipe_names
4 ['tok2vec', 'tagger', 'parser', 'attribute_ruler', 'lemmatizer', 'ner']

```

11828.py hosted with ❤️ by GitHub

[view raw](#)

When a text is processed with `nlp`, it is first tokenized and passed down to each pipeline from the above list, which, in turn, modifies and returns the Doc object with new information. Here is how it is illustrated in the Spacy docs:

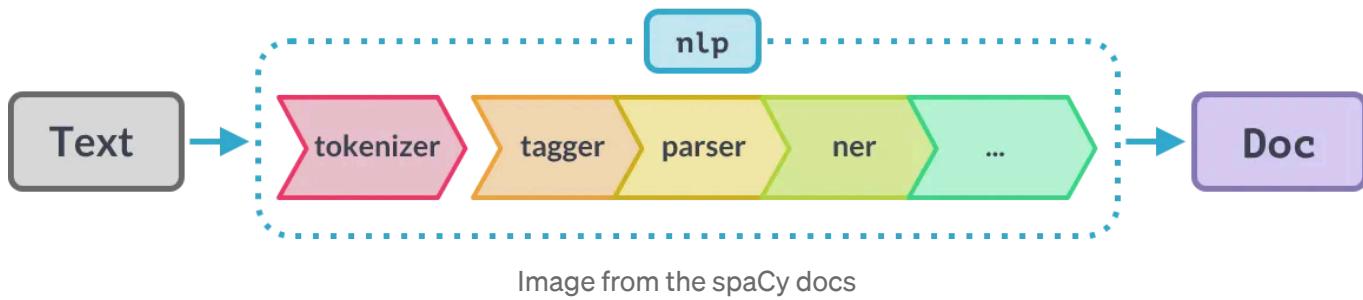


Image from the spaCy docs

But, spaCy doesn't have a custom pipeline for any NLP problem you might face in the real world. For example, you may want to perform your preprocessing steps before the text is passed to other pipelines or write callbacks to extract different types of information between pipelines.

For such cases, you should learn how to write custom pipeline functions and add them to the `nlp` object, so they are automatically run when you call `nlp` on text.

Here is a basic overview of how to do this:

```

1  from spacy.language import Language
2
3  @Language.component("your_component")
4  def your_component(doc):
5      # Do something on the doc
6      print(f"There are {len(doc)} tokens in this text.")
7
8  return doc

```

11829.py hosted with ❤️ by GitHub

[view raw](#)

You need the general `Language` class and decorate its `component` method over your function. The custom function must accept and return a single argument for the `doc` object. Above, we are defining a simple pipeline that prints out the length of the `doc` object.

Let's add it to the `nlp` object:

```

1  >>> nlp.add_pipe("your_component")
2  <function __main__.your_component(<doc>)
3
4  >>> nlp.pipe_names
5  ['tok2vec',
6  'tagger',
7  'parser',
8  'attribute_ruler',
9  'lemmatizer',
10 'ner',
11 'your_component']

```

11820.py hosted with ❤️ by GitHub

[view raw](#)

As you can see, the custom pipeline is added to the end. Now, we will call `nlp` on a sample text:

```
>>> doc = nlp("Bird dies, but you remember the flight.")
```

There are 9 tokens in this text.

Working as expected.

Now, let's do a more serious example. We will go back to "Alexander the Great" example and write a pipeline that adds the conqueror's name as an entity. And we want all this happen automatically without finding the entity outside the pipeline.

Here is the complete code:

```
1  from spacy.language import Language
2  from spacy.matcher import Matcher
3  from spacy.util import filter_spans
4
5  nlp = spacy.load("en_core_web_sm")
6
7  @Language.component("titled_person")
8  def titled_person(doc):
9      pattern = [
10          {"IS_ALPHA": True, "IS_TITLE": True},
11          {"IS_STOP": True},
12          {"IS_ALPHA": True, "IS_TITLE": True},
13      ]
14      # Create the matcher
15      matcher = Matcher(nlp.vocab)
16      # Add the pattern
17      matcher.add("TITLED_PERSON", [pattern])
18
19      matches = matcher(doc)
20      matched_spans = [Span(doc, start, end, label="PERSON") for _, start, end in matches]
21
22      # Filter the entities for potential overlap
23      filtered_matches = filter_spans(list(doc.ents) + matched_spans)
24      # Add the matched spans to doc's entities
25      doc.ents = filtered_matches
26
27      return doc
28
29  >>> nlp.add_pipe("titled_person")
30  <function __main__.titled_person(doc)>
```

11821.py hosted with ❤ by GitHub

[view raw](#)

We define the pattern, create a matcher object, add the matches to the entities, and return the doc. Let's test it:

```

1  txt = """Cleopatra wasn't actually Egyptian!
2      As far as historians can tell, Egypt's
3          famous femme fatal was actually Greek!.
4      She was a descendant of Alexander the Great's
5          Macedonian general Ptolemy"""
6
7  doc = nlp(txt)
8
9  >>> nlp.pipe_names
10 ['tok2vec',
11  'tagger',
12  'parser',
13  'attribute_ruler',
14  'lemmatizer',
15  'ner',
16  'titled_person']
17
18 >>> doc.ents
19 (Egyptian, Egypt, Greek, Alexander the Great, Ptolemy)

```

11833.py hosted with ❤️ by GitHub

[view raw](#)

It is working as expected.

So far, we have been appending custom pipelines to the end, but we can control this behavior. The `add_pipe` function has arguments to specify precisely where you want to insert the function:

```

1  nlp.add_pipe("titled_person", first=True) # Beginning
2  nlp.add_pipe("titled_person", after="parser") # After parser
3  nlp.add_pipe("titled_person", before="tagger") # Before POS tagger

```

11834.py hosted with ❤️ by GitHub

[view raw](#)

## Conclusion

Today, you have taken bold steps towards mastering the Scikit-learn of natural language processing. Armed with the article's knowledge, you can now roam freely across [spaCy's user guide](#), which is just as large and information-rich as Scikit-learn's.

Thank you for reading!

Loved this article and, let's face it, its bizarre writing style? Imagine having access to dozens more just like it, all written by a brilliant, charming, witty author (that's me, by the way :).

For only 4.99\$ membership, you will get access to not just my stories, but a treasure trove of knowledge from the best and brightest minds on Medium. And if you use [my referral link](#), you will earn my supernova of gratitude and a virtual high-five for supporting my work.

**Join Medium with my referral link — Bex T.**

Get exclusive access to all my  premium  content and all over Medium without limits. Support my work by buying me a...

[ibexorigin.medium.com](https://ibexorigin.medium.com)



[Artificial Intelligence](#)[Machine Learning](#)[Data Science](#)[Programming](#)[Editors Pick](#)

## Written by Bex T.

41K Followers · Writer for Towards Data Science

[Follow](#)

BEXGBoost | DataCamp Instructor | 🥇 Top 10 AI/ML Writer on Medium | Kaggle Master | <https://www.linkedin.com/in/bextuychiev/>

## More from Bex T. and Towards Data Science



Bex T. in Towards AI

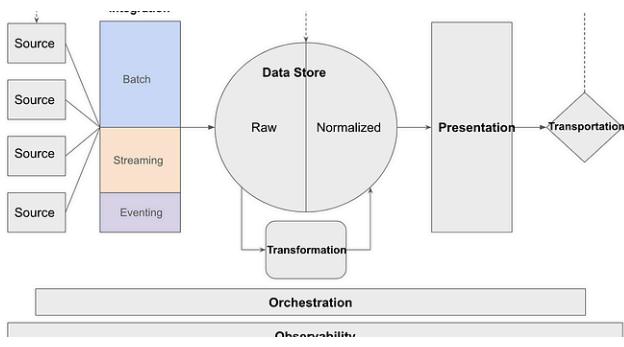
### 6 Pandas Mistakes That Silently Tell You Are a Rookie

No error messages—that's what makes them subtle

★ · 7 min read · Sep 4, 2023

2.1K 18

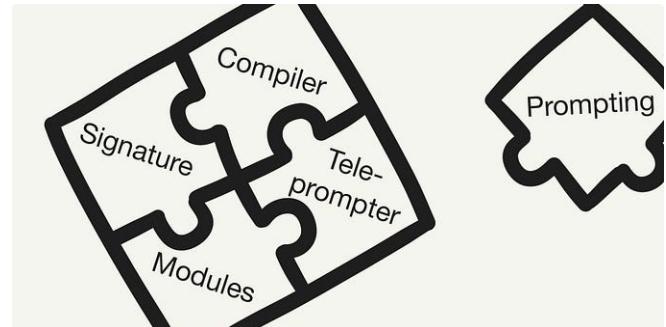
⋮



Dave Melillo in Towards Data Science

### Building a Data Platform in 2024

How to build a modern, scalable data platform to power your analytics and data science...



Leonie Monigatti in Towards Data Science

### Intro to DSPy: Goodbye Prompting, Hello Programming!

How the DSPy framework solves the fragility problem in LLM-based applications by...

★ · 13 min read · Feb 27, 2024

3.4K 10

⋮



Bex T. in Towards AI

### How to Boost Pandas Speed And Process 10M-row Datasets in...

Use Pandas the way it was intended to

9 min read · Feb 5, 2024

★ · 6 min read · Sep 7, 2023



2.4K



1.1K



...



...



5

[See all from Bex T.](#)[See all from Towards Data Science](#)

## Recommended from Medium

B-PER	O	O	B-LOC	I-LOC	O
John	lives	in	New	York	.

 Ahmet Münir Kocaman

### Mastering Named Entity Recognition with BERT: A...

Introduction

11 min read · Oct 5, 2023



...



126



...

10 min read · Dec 3, 2023

```
import spacy
nlp = spacy.load('textcat_singleLabel_model/model-best')

def compare_labels(fileid):
    text = reuters.raw(fileid)
    expected_cat = reuters.categories(fileid)

    doc = nlp(text)
    estimated_cats = sorted(doc.cats.items(), key=lambda i:float(i[1]), reverse=True)
    estimated_cat = estimated_cats[0]

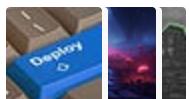
    print(f'{fileid} :: expected = {expected_cat} // estimated = {estimated_cat}'')
```

 Sebastian

### NLP with Spacy: Custom Text Classification Pipeline

Spacy is a powerful NLP library in which many NLP tasks like tokenization, stemming, part-

## Lists



## Predictive Modeling w/ Python

20 stories · 1013 saves



## Natural Language Processing

1305 stories · 793 saves



## Practical Guides to Machine Learning

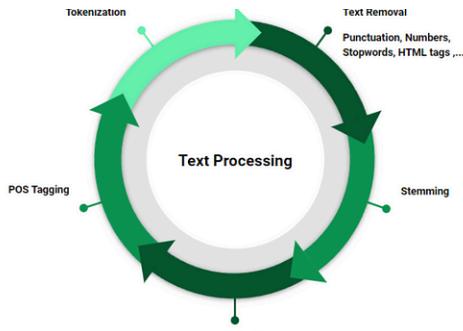
10 stories · 1219 saves



## ChatGPT prompts

47 stories · 1291 saves

 Christian Thrun PERSON started working on self-driv  
 in 2007 DATE , few people outside of the co



 Isidoro Grisanti

## Named Entity Recognition with LLMs—Extract Conversation...

This article aims to provide a comprehensive overview how to solve NER tasks based on...

5 min read · Oct 20, 2023



47



...



 Johni Douglas Marangon

 Aysel Aydin

## 1—Text Preprocessing Techniques for NLP

In this article, we will cover the following topics:

4 min read · Oct 4, 2023



259



...



 Yuan An, PhD

## Building a custom Named Entity Recognition model using spaCy ...

In today's post, we will learn how to train a NER. In the previous post we saw the...

4 min read · Nov 21, 2023

👏 60

💬 1



...

## Use bert-base-NER in Hugging Face for Named Entity Recognition

Named Entity Recognition (NER) is a subtask of information extraction that classifies...

⭐ · 3 min read · Sep 23, 2023

👏 17

💬



...

See more recommendations