




Doc

CLASS

[</> SOURCE](#)

A container for accessing linguistic annotations.

A `Doc` is a sequence of `Token`  objects. Access sentences and named entities, export annotations to numpy arrays, losslessly serialize to compressed binary strings. The `Doc` object holds an array of `TokenC`  structs. The Python-level `Token` and `Span`  objects are views of this array, i.e. they don't own the data themselves.

Doc.__init__

METHOD

Construct a `Doc` object. The most common way to get a `Doc` object is via the `nlp` object.

NAME	DESCRIPTION
vocab	A storage container for lexical types. TYPE: <code>Vocab</code>
words	A list of strings or integer hash values to add to the document as words. TYPE: <code>Optional[List[Union[str,int]]]</code>
spaces	A list of boolean values indicating whether each word has a subsequent space. Must have the same length as <code>words</code> , if specified. Defaults to a sequence of <code>True</code> . TYPE: <code>Optional[List[bool]]</code>
KEYWORD-ONLY	
user_data	Optional extra data to attach to the Doc. TYPE: <code>Dict</code>
tags V3.0 ?	A list of strings, of the same length as <code>words</code> , to assign as <code>token.tag</code> for each word. Defaults to <code>None</code> . TYPE: <code>Optional[List[str]]</code>
pos V3.0 ?	A list of strings, of the same length as <code>words</code> , to assign as <code>token.pos</code> for each word. Defaults to <code>None</code> . TYPE: <code>Optional[List[str]]</code>
morphs V3.0 ?	A list of strings, of the same length as <code>words</code> , to assign as <code>token.morph</code> for each word. Defaults to <code>None</code> . TYPE: <code>Optional[List[str]]</code>
lemmas V3.0 ?	A list of strings, of the same length as <code>words</code> , to assign as <code>token.lemma</code> for each word. Defaults to <code>None</code> . TYPE: <code>Optional[List[str]]</code>
heads V3.0 ?	A list of values, of the same length as <code>words</code> , to assign as the head for each word. Head indices are the absolute position of the head in the <code>Doc</code> . Defaults to <code>None</code> . TYPE: <code>Optional[List[int]]</code>
deps V3.0 ?	A list of strings, of the same length as <code>words</code> , to assign as <code>token.dep</code> for each word. Defaults to <code>None</code> . TYPE: <code>Optional[List[str]]</code>
sent_starts V3.0 ?	A list of values, of the same length as <code>words</code> , to assign as <code>token.is_sent_start</code> . Will be overridden by <code>heads</code> if <code>heads</code> is provided. Defaults to <code>None</code> . TYPE: <code>Optional[List[Union[bool, int, None]]]</code>
ents V3.0 ?	A list of strings, of the same length of <code>words</code> , to assign the token-based IOB tag. Defaults to <code>None</code> .

Defaults to `None`.

TYPE: `Optional[List[str]]`

Doc.__getitem__ METHOD

Get a `Token` object at position `i`, where `i` is an integer. Negative indexing is supported, and follows the usual Python semantics, i.e. `doc[-2]` is `doc[len(doc) - 2]`.


NAME	DESCRIPTION
<code>i</code>	The index of the token. TYPE: <code>int</code>
RETURNS	The token at <code>doc[i]</code> . TYPE: <code>Token</code>

Get a `Span` object, starting at position `start` (token index) and ending at position `end` (token index). For instance, `doc[2:5]` produces a span consisting of tokens 2, 3 and 4. Stepped slices (e.g. `doc[start : end : step]`) are not supported, as `Span` objects must be contiguous (cannot have gaps). You can use negative indices and open-ended ranges, which have their normal Python semantics.

NAME	DESCRIPTION
<code>start_end</code>	The slice of the document to get. TYPE: <code>Tuple[int, int]</code>
RETURNS	The span at <code>doc[start:end]</code> . TYPE: <code>Span</code>

Doc.__iter__ METHOD

Iterate over `Token` objects, from which the annotations can be easily accessed.

This is the main way of accessing `Token`  objects, which are the main way annotations are accessed from Python. If faster-than-Python speeds are required, you can instead access the annotations as a numpy array, or access the underlying C data directly from Cython.

NAME	DESCRIPTION
------	-------------

YIELDS	A <code>Token</code> object.
	TYPE: <code>Token</code>

Doc.__len__ METHOD

Get the number of tokens in the document.

NAME	DESCRIPTION
------	-------------

RETURNS	The number of tokens in the document.
	TYPE: <code>int</code>

Doc.set_extension CLASSMETHOD

Define a custom attribute on the `Doc` which becomes available via `Doc._`. For details, see the documentation on [custom attributes](#).

NAME	DESCRIPTION
name	Name of the attribute to set by the extension. For example, "my_attr" will be available as <code>doc._.my_attr</code> . TYPE: str
default	Optional default value of the attribute if no getter or method is defined. TYPE: Optional[Any]
method	Set a custom method on the object, for example <code>doc._.compare(other_doc)</code> . TYPE: Optional[Callable[[Doc, ...], Any]]
getter	Getter function that takes the object and returns an attribute value. Is called when the user accesses the <code>._</code> attribute. TYPE: Optional[Callable[[Doc], Any]]
setter	Setter function that takes the <code>Doc</code> and a value, and modifies the object. Is called when the user writes to the <code>Doc._</code> attribute. TYPE: Optional[Callable[[Doc, Any], None]]
force	Force overwriting existing attribute. TYPE: bool

Doc.get_extension CLASSMETHOD

Look up a previously registered extension by name. Returns a 4-tuple (default, method, getter, setter) if the extension is registered. Raises a `KeyError` otherwise.

NAME	DESCRIPTION
<code>name</code>	Name of the extension. TYPE: <code>str</code>
RETURNS	A <code>(default, method, getter, setter)</code> tuple of the extension. TYPE: <code>Tuple[Optional[Any], Optional[Callable], Optional[Callable], Optional[Callable]]</code>

Doc.has_extension CLASSMETHOD

Check whether an extension has been registered on the `Doc` class.

NAME	DESCRIPTION
<code>name</code>	Name of the extension to check. TYPE: <code>str</code>
RETURNS	Whether the extension has been registered. TYPE: <code>bool</code>

Doc.remove_extension CLASSMETHOD

Remove a previously registered extension.

NAME	DESCRIPTION
<code>name</code>	Name of the extension. TYPE: <code>str</code>
RETURNS	A <code>(default, method, getter, setter)</code> tuple of the removed extension. TYPE: <code>Tuple[Optional[Any], Optional[Callable], Optional[Callable], Optional[Callable]]</code>

Doc.char_span METHOD

Create a `Span` object from the slice `doc.text[start_idx:end_idx]`. Returns `None` if the character indices don't map to a valid span using the default alignment mode `"strict"`.

NAME	DESCRIPTION
<code>start</code>	The index of the first character of the span. TYPE: <code>int</code>
<code>end</code>	The index of the last character after the span. TYPE: <code>int</code>
<code>label</code>	A label to attach to the span, e.g. for named entities. TYPE: <code>Union[int, str]</code>
<code>kb_id</code>	An ID from a knowledge base to capture the meaning of a named entity. TYPE: <code>Union[int, str]</code>
<code>vector</code>	A meaning representation of the span. TYPE: <code>numpy.ndarray[ndim=1, dtype=float32]</code>
<code>alignment_mode</code>	How character indices snap to token boundaries. Options: <code>"strict"</code> (no snapping), <code>"contract"</code> (span of all tokens completely within the character span), <code>"expand"</code> (span of all tokens at least partially covered by the character span). Defaults to <code>"strict"</code> . TYPE: <code>str</code>
<code>span_id</code> V3.3.1 ?	An identifier to associate with the span. TYPE: <code>Union[int, str]</code>
RETURNS	The newly constructed object or <code>None</code> . TYPE: <code>Optional[Span]</code>

Doc.set_ents METHOD V3.0 ?

Set the named entities in the document.

NAME	DESCRIPTION
<code>entities</code>	Spans with labels to set as entities. TYPE: <code>List[Span]</code>
KEYWORD-ONLY	
<code>blocked</code>	Spans to set as “blocked” (never an entity) for spacy’s built-in NER component. Other components may ignore this setting. TYPE: <code>Optional[List[Span]]</code>
<code>missing</code>	Spans with missing/unknown entity information. TYPE: <code>Optional[List[Span]]</code>
<code>outside</code>	Spans outside of entities (O in IOB). TYPE: <code>Optional[List[Span]]</code>
<code>default</code>	How to set entity annotation for tokens outside of any provided spans. Options: <code>"blocked"</code> , <code>"missing"</code> , <code>"outside"</code> and <code>"unmodified"</code> (preserve current state). Defaults to <code>"outside"</code> . TYPE: <code>str</code>

Doc.similarity

METHOD

NEEDS MODEL ?

Make a semantic similarity estimate. The default estimate is cosine similarity using an average of word vectors.

NAME	DESCRIPTION
<code>other</code>	The object to compare with. By default, accepts <code>Doc</code> , <code>Span</code> , <code>Token</code> and <code>Lexeme</code> objects. TYPE: <code>Union[Doc, Span, Token, Lexeme]</code>
RETURNS	A scalar similarity score. Higher is more similar. TYPE: <code>float</code>

Doc.count_by METHOD

Count the frequencies of a given attribute. Produces a dict of `{attr (int): count (ints)}` frequencies, keyed by the values of the given attribute ID.


NAME	DESCRIPTION
<code>attr_id</code>	The attribute ID. TYPE: <code>int</code>
RETURNS	A dictionary mapping attributes to integer counts. TYPE: <code>Dict[int, int]</code>

Doc.get_lca_matrix METHOD

Calculates the lowest common ancestor matrix for a given `Doc`. Returns LCA matrix containing the integer index of the ancestor, or `-1` if no common ancestor is found, e.g. if span excludes a necessary ancestor.

NAME	DESCRIPTION
RETURNS	The lowest common ancestor matrix of the <code>Doc</code> . TYPE: <code>numpy.ndarray[ndim=2, dtype=int32]</code>

Doc.has_annotation METHOD

Check whether the doc contains annotation on a `Token attribute` .

NAME	DESCRIPTION
<code>attr</code>	The attribute string name or int ID. TYPE: <code>Union[int, str]</code>
KEYWORD-ONLY <code>require_complete</code>	Whether to check that the attribute is set on every token in the doc. Defaults to <code>False</code> . TYPE: <code>bool</code>
RETURNS	Whether specified annotation is present in the doc. TYPE: <code>bool</code>

Doc.to_array **METHOD**

Export given token attributes to a numpy `ndarray`. If `attr_ids` is a sequence of `M` attributes, the output array will be of shape `(N, M)`, where `N` is the length of the `Doc` (in tokens). If `attr_ids` is a single attribute, the output shape will be `(N,)`. You can specify attributes by integer ID (e.g. `spacy.attrs.LEMMA`) or string name (e.g. "LEMMA" or "lemma"). The values will be 64-bit integers.

Returns a 2D array with one row per token and one column per attribute (when `attr_ids` is a list), or as a 1D numpy array, with one item per attribute (when `attr_ids` is a single value).

NAME	DESCRIPTION
<code>attr_ids</code>	A list of attributes (int IDs or string names) or a single attribute (int ID or string name). TYPE: <code>Union[int, str, List[Union[int, str]]]</code>
RETURNS	The exported attributes as a numpy array. TYPE: <code>Union[numpy.ndarray[ndim=2, dtype=uint64], numpy.ndarray[ndim=1, dtype=uint64]]</code>


Doc.from_array METHOD

Load attributes from a numpy array. Write to a `Doc` object, from an `(M, N)` array of attributes.

NAME	DESCRIPTION
<code>attrs</code>	A list of attribute ID ints. TYPE: <code>List[int]</code>
<code>array</code>	The attribute values to load. TYPE: <code>numpy.ndarray[ndim=2, dtype=int32]</code>
<code>exclude</code>	String names of serialization fields to exclude. TYPE: <code>Iterable[str]</code>
RETURNS	The <code>Doc</code> itself. TYPE: <code>Doc</code>

Doc.from_docs STATICMETHOD V3.0 ?

Concatenate multiple `Doc` objects to form a new one. Raises an error if the `Doc` objects do not all share the same `Vocab`.

NAME	DESCRIPTION
<code>docs</code>	<p>A list of <code>Doc</code> objects.</p> <p>TYPE: <code>List[Doc]</code></p>
<code>ensure_whitespace</code>	<p>Insert a space between two adjacent docs whenever the first doc does not end in whitespace.</p> <p>TYPE: <code>bool</code></p>
<code>attrs</code>	<p>Optional list of attribute ID ints or attribute name strings.</p> <p>TYPE: <code>Optional[List[Union[str, int]]]</code></p>
KEYWORD-ONLY	
<code>exclude</code> V3.3 	<p>String names of Doc attributes to exclude. Supported: <code>spans</code>, <code>tensor</code>, <code>user_data</code>.</p> <p>TYPE: <code>Iterable[str]</code></p>
RETURNS	<p>The new <code>Doc</code> object that is containing the other docs or <code>None</code>, if <code>docs</code> is empty or <code>None</code>.</p> <p>TYPE: <code>Optional[Doc]</code></p>

Doc.to_disk **METHOD**

Save the current state to a directory.

NAME	DESCRIPTION
<code>path</code>	<p>A path to a directory, which will be created if it doesn't exist. Paths may be either strings or <code>Path</code>-like objects.</p> <p>TYPE: <code>Union[str, Path]</code></p>
KEYWORD-ONLY	
<code>exclude</code>	<p>String names of <u>serialization fields</u> to exclude.</p> <p>TYPE: <code>Iterable[str]</code></p>

Doc.from_disk METHOD

Loads state from a directory. Modifies the object in place and returns it.

NAME	DESCRIPTION
<code>path</code>	A path to a directory. Paths may be either strings or <code>Path</code> -like objects. TYPE: <code>Union[str, Path]</code>
KEYWORD-ONLY	
<code>exclude</code>	String names of serialization fields to exclude. TYPE: <code>Iterable[str]</code>
RETURNS	The modified <code>Doc</code> object. TYPE: <code>Doc</code>

Doc.to_bytes METHOD

Serialize, i.e. export the document contents to a binary string.


NAME	DESCRIPTION
KEYWORD-ONLY	
<code>exclude</code>	String names of serialization fields to exclude. TYPE: <code>Iterable[str]</code>
RETURNS	A losslessly serialized copy of the <code>Doc</code> , including all annotations. TYPE: <code>bytes</code>

Doc.from_bytes METHOD

Deserialize, i.e. import the document contents from a binary string.


NAME	DESCRIPTION
<code>data</code>	The string to load from. TYPE: <code>bytes</code>
KEYWORD-ONLY	
<code>exclude</code>	String names of serialization fields to exclude. TYPE: <code>Iterable[str]</code>
RETURNS	The <code>Doc</code> object.
	TYPE: <code>Doc</code>

Doc.to_json METHOD

Serializes a document to JSON. Note that this is format differs from the deprecated `JSON training format`  .

NAME	DESCRIPTION
<code>underscore</code>	Optional list of string names of custom <code>Doc</code> attributes. Attribute values need to be JSON-serializable. Values will be added to an <code>"_"</code> key in the data, e.g. <code>"_": {"foo": "bar"}</code> . TYPE: <code>Optional[List[str]]</code>
RETURNS	The data in JSON format.
	TYPE: <code>Dict[str, Any]</code>

Doc.from_json METHOD V3.3.1

Deserializes a document from JSON, i.e. generates a document from the provided JSON data as generated by `Doc.to_json()`  .

NAME	DESCRIPTION
<code>doc_json</code>	The Doc data in JSON format from <code>Doc.to_json</code> . TYPE: Dict[str, Any]
KEYWORD-ONLY	
<code>validate</code>	Whether to validate the JSON input against the expected schema for detailed debugging. Defaults to <code>False</code> . TYPE: bool
RETURNS	A <code>Doc</code> corresponding to the provided JSON. TYPE: <code>Doc</code>

Doc.retokenize CONTEXTMANAGER

Context manager to handle retokenization of the `Doc` . Modifications to the `Doc` 's tokenization are stored, and then made all at once when the context manager exits. This is much more efficient, and less error-prone. All views of the `Doc` (`Span` and `Token`) created before the retokenization are invalidated, although they may accidentally continue to work.

NAME	DESCRIPTION
RETURNS	The retokenizer. TYPE: Retokenizer

Retokenizer.merge METHOD

Mark a span for merging. The `attrs` will be applied to the resulting token (if they're context-dependent token attributes like `LEMMA` or `DEP`) or to the underlying lexeme (if they're context-independent lexical attributes like `LOWER` or `IS_STOP`). Writable custom extension attributes can be provided using the `"_"` key and specifying a dictionary that maps attribute names to values.

NAME	DESCRIPTION
<code>span</code>	The span to merge. TYPE: <code>Span</code>
<code>attrs</code>	Attributes to set on the merged token. TYPE: <code>Dict[Union[str, int], Any]</code>

Retokenizer.split METHOD

Mark a token for splitting, into the specified `orths`. The `heads` are required to specify how the new subtokens should be integrated into the dependency tree. The list of per-token heads can either be a token in the original document, e.g. `doc[2]`, or a tuple consisting of the token in the original document and its subtoken index. For example, `(doc[3], 1)` will attach the subtoken to the second subtoken of `doc[3]`.

This mechanism allows attaching subtokens to other newly created subtokens, without having to keep track of the changing token indices. If the specified head token will be split within the retokenizer block and no subtoken index is specified, it will default to `0`. Attributes to set on subtokens can be provided as a list of values. They'll be applied to the resulting token (if they're context-dependent token attributes like `LEMMA` or `DEP`) or to the underlying lexeme (if they're context-independent lexical attributes like `LOWER` or `IS_STOP`).

NAME	DESCRIPTION
<code>token</code>	The token to split. TYPE: <code>Token</code>
<code>orths</code>	The verbatim text of the split tokens. Needs to match the text of the original token. TYPE: <code>List[str]</code>
<code>heads</code>	List of <code>token</code> or <code>(token, subtoken)</code> tuples specifying the tokens to attach the newly split subtokens to. TYPE: <code>List[Union[Token, Tuple[Token, int]]]</code>
<code>attrs</code>	Attributes to set on all split tokens. Attribute names mapped to list of per-token attribute values. TYPE: <code>Dict[Union[str, int], List[Any]]</code>

Doc.ents

PROPERTY

NEEDS MODEL ?

The named entities in the document. Returns a tuple of named entity `Span` objects, if the entity recognizer has been applied.

NAME	DESCRIPTION
------	-------------

RETURNS	Entities in the document, one <code>Span</code> per entity.
---------	---

TYPE: `Tuple[Span]`

Doc.spans

PROPERTY

A dictionary of named span groups, to store and access additional span annotations. You can write to it by assigning a list of `Span` objects or a `SpanGroup` to a given key.

NAME	DESCRIPTION
------	-------------

RETURNS	The span groups assigned to the document.
---------	---

TYPE: `Dict[str, SpanGroup]`

Doc.cats

PROPERTY

NEEDS MODEL ?

Maps a label to a score for categories applied to the document. Typically set by the `TextCategorizer`.

NAME	DESCRIPTION
------	-------------

RETURNS	The text categories mapped to scores.
	TYPE: Dict[str, float]

Doc.noun_chunks

PROPERTY

NEEDS MODEL ?

Iterate over the base noun phrases in the document. Yields base noun-phrase `Span` objects, if the document has been syntactically parsed. A base noun phrase, or “NP chunk”, is a noun phrase that does not permit other NPs to be nested within it – so no NP-level coordination, no prepositional phrases, and no relative clauses.

To customize the noun chunk iterator in a loaded pipeline, modify `nlp.vocab.get_noun_chunks` . If the `noun_chunk` [syntax iterator](#) has not been implemented for the given language, a `NotImplementedError` is raised.

NAME	DESCRIPTION
------	-------------

YIELDS	Noun chunks in the document.
	TYPE: <code>Span</code>

Doc.sents

PROPERTY

NEEDS MODEL ?

Iterate over the sentences in the document. Sentence spans have no label.

This property is only available when [sentence boundaries](#) have been set on the document by the `parser` , `sender` , `sentencizer` or some custom function. It will raise an error otherwise.

NAME	DESCRIPTION
------	-------------

YIELDS	Sentences in the document.
	TYPE: <code>Span</code>

Doc.has_vector

PROPERTY

NEEDS MODEL ?

A boolean value indicating whether a word vector is associated with the object.

NAME	DESCRIPTION
------	-------------

RETURNS	Whether the document has a vector data attached.
	TYPE: <code>bool</code>

Doc.vector

PROPERTY

NEEDS MODEL ?

A real-valued meaning representation. Defaults to an average of the token vectors.

NAME	DESCRIPTION
------	-------------

RETURNS	A 1-dimensional array representing the document's vector.
	TYPE: <code>numpy.ndarray[ndim=1, dtype=float32]</code>

Doc.vector_norm

PROPERTY

NEEDS MODEL ?

The L2 norm of the document's vector representation.

NAME

DESCRIPTION

RETURNS

The L2 norm of the vector representation.

TYPE: float

Attributes

NAME	DESCRIPTION
<code>text</code>	A string representation of the document text. TYPE: <code>str</code>
<code>text_with_ws</code>	An alias of <code>Doc.text</code> , provided for duck-type compatibility with <code>Span</code> and <code>Token</code> . TYPE: <code>str</code>
<code>mem</code>	The document's local memory heap, for all C data it owns. TYPE: <code>cymem.Pool</code>
<code>vocab</code>	The store of lexical types. TYPE: <code>Vocab</code>
<code>tensor</code>	Container for dense vector representations. TYPE: <code>numpy.ndarray</code>
<code>user_data</code>	A generic storage area, for user custom data. TYPE: <code>Dict[str, Any]</code>
<code>lang</code>	Language of the document's vocabulary. TYPE: <code>int</code>
<code>lang_</code>	Language of the document's vocabulary. TYPE: <code>str</code>
<code>sentiment</code>	The document's positivity/negativity score, if available. TYPE: <code>float</code>
<code>user_hooks</code>	A dictionary that allows customization of the <code>Doc</code> 's properties. TYPE: <code>Dict[str, Callable]</code>
<code>user_token_hooks</code>	A dictionary that allows customization of properties of <code>Token</code> children. TYPE: <code>Dict[str, Callable]</code>
<code>user_span_hooks</code>	A dictionary that allows customization of properties of <code>Span</code> children. TYPE: <code>Dict[str, Callable]</code>
<code>has_unknown_spaces</code>	Whether the document was constructed without known spacing between tokens (typically when created from gold tokenization). TYPE: <code>bool</code>
<code>_</code>	User space for adding custom attribute extensions .

Serialization fields

During serialization, spaCy will export several data fields used to restore different aspects of the object. If needed, you can exclude them from serialization by passing in the string names via the `exclude` argument.

NAME	DESCRIPTION
<code>text</code>	The value of the <code>Doc.text</code> attribute.
<code>sentiment</code>	The value of the <code>Doc.sentiment</code> attribute.
<code>tensor</code>	The value of the <code>Doc.tensor</code> attribute.
<code>user_data</code>	The value of the <code>Doc.user_data</code> dictionary.
<code>user_data_keys</code>	The keys of the <code>Doc.user_data</code> dictionary.
<code>user_data_values</code>	The values of the <code>Doc.user_data</code> dictionary.

[</> SUGGEST EDITS](#)