

---

# antu Documentation

*Release 0.0.1*

**AntNLP**

**May 24, 2020**



---

## Contents:

---

<b>1</b>	<b>antu.io package</b>	<b>3</b>
1.1	Subpackages . . . . .	3
1.2	Submodules . . . . .	9
1.3	antu.io.instance module . . . . .	9
1.4	antu.io.vocabulary module . . . . .	11
1.5	Module contents . . . . .	13
<b>2</b>	<b>antu.nn package</b>	<b>15</b>
2.1	Subpackages . . . . .	15
2.2	Module contents . . . . .	15
<b>3</b>	<b>Indices and tables</b>	<b>17</b>
	<b>Python Module Index</b>	<b>19</b>
	<b>Index</b>	<b>21</b>



Universal data IO and neural network modules in NLP tasks.

- \*data IO is an universal module in Natural Language Processing system and not based on any framework (like TensorFlow, PyTorch, MXNet, Dynet. . .).
- \*neural network module contains the neural network structures commonly used in NLP tasks. We want to design commonly used structures for each neural network framework. We will continue to develop this module.



## 1.1 Subpackages

### 1.1.1 antu.io.dataset\_readers package

#### Submodules

#### antu.io.dataset\_readers.dataset\_reader module

**class** antu.io.dataset\_readers.dataset\_reader.**DatasetReader**  
Bases: `object`

#### Methods

<b>input_to_instance</b>	
<b>read</b>	

**input\_to\_instance** (*inputs: str*) → antu.io.instance.Instance

**read** (*file\_path: str*) → List[antu.io.instance.Instance]

#### Module contents

### 1.1.2 antu.io.datasets package

#### Submodules

## antu.io.datasets.dataset module

**class** antu.io.datasets.dataset.**Dataset**  
 Bases: `object`

### Methods

<b>build_dataset</b>	
----------------------	--

```
build_dataset()
datasets = {}
vocabulary_set = {}
```

### Module contents

## 1.1.3 antu.io.fields package

### Submodules

#### antu.io.fields.field module

**class** antu.io.fields.field.**Field**  
 Bases: `object`

A `Field` is an ingredient of a data instance. In most NLP tasks, `Field` stores data of string types. It contains one or more indexers that map string data to the corresponding index. Data instances are collections of fields.

### Methods

<code>count_vocab_items(counter, Dict[str, int])</code>	We count the number of strings if the string needs to be mapped to one or more integers.
<code>index(vocab)</code>	Gets one or more index mappings for each element in the <code>Field</code> .

**count\_vocab\_items** (*counter: Dict[str, Dict[str, int]]*) → None  
 We count the number of strings if the string needs to be mapped to one or more integers. You can pass directly if there is no string that needs to be mapped.

#### Parameters

**counter** [Dict[str, Dict[str, int]]]

“counter“ is used to count the number of each item. The first key represents the namespace of the vocabulary, and the second key represents the string of the item.

**index** (*vocab: antu.io.vocabulary.Vocabulary*) → None  
 Gets one or more index mappings for each element in the `Field`.

#### Parameters



**vocab** [Vocabulary]

“vocab“ is used to get the index of each item.

### antu.io.fields.index\_field module

**class** antu.io.fields.index\_field.**IndexField** (*name: str, tokens: List[str]*)

Bases: *antu.io.fields.field.Field*

A IndexField is an integer field, and we can use it to store data ID.

#### Parameters

**name** [str] Field name. This is necessary and must be unique (not the same as other field names).

**tokens** [List[str]] Field content that contains a list of string.

#### Methods

<i>count_vocab_items</i> (counters, Dict[str, int])	IndexField doesn't need index operation.
<i>index</i> (vocab)	IndexField doesn't need index operation.

**count\_vocab\_items** (*counters: Dict[str, Dict[str, int]]*) → None  
IndexField doesn't need index operation.

**index** (*vocab: antu.io.vocabulary.Vocabulary*) → None  
IndexField doesn't need index operation.

### antu.io.fields.sequence\_label\_field module

**class** antu.io.fields.sequence\_label\_field.**SequenceLabelField** (*name: str, tokens: List[str], indexers:*

*List[antu.io.token\_indexers.token\_indexer.T*

Bases: *antu.io.fields.field.Field*

#### Methods

<i>count_vocab_items</i> (counters, Dict[str, int])	We count the number of strings if the string needs to be mapped to one or more integers.
<i>index</i> (vocab)	Gets one or more index mappings for each element in the Field.

**count\_vocab\_items** (*counters: Dict[str, Dict[str, int]]*) → None  
We count the number of strings if the string needs to be mapped to one or more integers. You can pass directly if there is no string that needs to be mapped.

#### Parameters

**counter** [Dict[str, Dict[str, int]]]

“counter“ is used to count the number of each item. The first key

represents the namespace of the vocabulary, and the second key represents

the string of the item.

**index** (*vocab*: *antu.io.vocabulary.Vocabulary*) → None  
 Gets one or more index mappings for each element in the Field.

**Parameters**

**vocab** [*Vocabulary*]  
 “vocab“ is used to get the index of each item.

**antu.io.fields.text\_field module**

**class** *antu.io.fields.text\_field.TextField* (*name*: *str*, *tokens*: *List[str]*, *indexers*:  
*List[antu.io.token\_indexers.token\_indexer.TokenIndexer]*  
 = [])

Bases: *antu.io.fields.field.Field*

A *TextField* is a data field that is commonly used in NLP tasks, and we can use it to store text sequences such as sentences, paragraphs, POS tags, and so on.

**Parameters**

**name** [*str*] Field name. This is necessary and must be unique (not the same as other field names).  
**tokens** [*List[str]*] Field content that contains a list of string.  
**indexers** [*List[TokenIndexer]*], optional (default=“list(“)) Indexer list that defines the vocabularies associated with the field.

**Methods**

---

<i>count_vocab_items</i> (counters, Dict[str, int])	We count the number of strings if the string needs to be counted to some
<i>index</i> (vocab)	Gets one or more index mappings for each element in the Field.

---

**count\_vocab\_items** (*counters*: *Dict[str, Dict[str, int]]*) → None

**We count the number of strings if the string needs to be counted to some** counters. You can pass directly if there is no string that needs to be counted.

**Parameters**

**counters** [*Dict[str, Dict[str, int]]*] Element statistics for datasets. if field indexers indicate that this field is related to some counters, we use field content to update the counters.

**index** (*vocab*: *antu.io.vocabulary.Vocabulary*) → None  
 Gets one or more index mappings for each element in the Field.

**Parameters**

**vocab** [*Vocabulary*] vocab is used to get the index of each item.

## Module contents

### 1.1.4 antu.io.token\_indexers package

#### Submodules

#### antu.io.token\_indexers.char\_token\_indexer module

```
class antu.io.token_indexers.char_token_indexer.CharTokenIndexer (related_vocabs:
    List[str],
    transform:
    Callable[[str],
    str] = <function CharTokenIndexer.<lambda>>>)
```

Bases: *antu.io.token\_indexers.token\_indexer.TokenIndexer*

A CharTokenIndexer determines how string token get represented as arrays of list of character indices in a model.

#### Parameters

**related\_vocabs** [*List[str]*] Which vocabularies are related to the indexer.

**transform** [*Callable[[str, ], str]*, optional (default=`'lambda x:x'`)] What changes need to be made to the token when counting or indexing. Commonly used are lowercase transformation functions.

#### Methods

<i>count_vocab_items</i> (token, counters, Dict[str, ...])	Each character in the token is counted directly as an element.
<i>tokens_to_indices</i> (tokens, vocab)	Takes a list of tokens and converts them to one or more sets of indices.

**count\_vocab\_items** (*token: str, counters: Dict[str, Dict[str, int]]*) → None  
Each character in the token is counted directly as an element.

#### Parameters

**counter** [*Dict[str, Dict[str, int]]*] We count the number of strings if the string needs to be counted to some counters.

**tokens\_to\_indices** (*tokens: List[str], vocab: antu.io.vocabulary.Vocabulary*) → *Dict[str, List[List[int]]]*

Takes a list of tokens and converts them to one or more sets of indices. During the indexing process, each token item corresponds to a list of index in the vocabulary.

#### Parameters

**vocab** [*Vocabulary*] vocab is used to get the index of each item.

## antu.io.token\_indexers.single\_id\_token\_indexer module

```
class antu.io.token_indexers.single_id_token_indexer.SingleIdTokenIndexer (related_vocabs:
                                                                    List[str],
                                                                    trans-
                                                                    form:
                                                                    Callable[[str],
                                                                    str]
                                                                    =
                                                                    <func-
                                                                    tion
                                                                    Sin-
                                                                    gleI-
                                                                    d-
                                                                    To-
                                                                    kenIn-
                                                                    dexer.<lambda>>)
```

Bases: *antu.io.token\_indexers.token\_indexer.TokenIndexer*

A SingleIdTokenIndexer determines how string token get represented as arrays of single id indices in a model.

### Parameters

**related\_vocabs** [List[str]] Which vocabularies are related to the indexer.

**transform** [Callable[[str], str], optional (default="lambda x:x")] What changes need to be made to the token when counting or indexing. Commonly used are lowercase transformation functions.

### Methods

---

<i>count_vocab_items</i> (token, counters, Dict[str, ...])	The token is counted directly as an element.
<i>tokens_to_indices</i> (tokens, vocab)	Takes a list of tokens and converts them to one or more sets of indices.

---

**count\_vocab\_items** (*token: str, counters: Dict[str, Dict[str, int]]*) → None  
The token is counted directly as an element.

### Parameters

**counter** [Dict[str, Dict[str, int]]] We count the number of strings if the string needs to be counted to some counters.

**tokens\_to\_indices** (*tokens: List[str], vocab: antu.io.vocabulary.Vocabulary*) → Dict[str, List[int]]  
Takes a list of tokens and converts them to one or more sets of indices. During the indexing process, each item corresponds to an index in the vocabulary.

### Parameters

**vocab** [Vocabulary] vocab is used to get the index of each item.

### Returns

**res** [Dict[str, List[int]]] if the token and index list is [w1:5, w2:3, w3:0], the result will be {'vocab\_name': [5, 3, 0]}

## antu.io.token\_indexers.token\_indexer module

**class** antu.io.token\_indexers.token\_indexer.**TokenIndexer**

Bases: `object`

A `TokenIndexer` determines how string tokens get represented as arrays of indices in a model.

### Methods

<code>count_vocab_items(token, counter, Dict[str, ...])</code>	Defines how each token in the field is counted.
<code>tokens_to_indices(tokens, vocab)</code>	Takes a list of tokens and converts them to one or more sets of indices.

**count\_vocab\_items** (*token: str, counter: Dict[str, Dict[str, int]]*) → None

Defines how each token in the field is counted. In most cases, just use the string as a key. However, for character-level `TokenIndexer`, you need to traverse each character in the string.

#### Parameters

**counter** [Dict[str, Dict[str, int]]] We count the number of strings if the string needs to be counted to some counters.

**tokens\_to\_indices** (*tokens: List[str], vocab: antu.io.vocabulary.Vocabulary*) → Dict[str, Indices]

Takes a list of tokens and converts them to one or more sets of indices. This could be just an ID for each token from the vocabulary.

#### Parameters

**vocab** [Vocabulary] `vocab` is used to get the index of each item.

## Module contents

### 1.2 Submodules

### 1.3 antu.io.instance module

**class** antu.io.instance.**Instance** (*fields: List[antu.io.fields.field.Field] = None*)

Bases: `collections.abc.Mapping`, `typing.Generic`

An `Instance` is a collection (list) of multiple data fields.

#### Parameters

**fields** [List[Field], optional (default="None")] A list of multiple data fields.

### Methods

<code>add_field(field)</code>	Add the field to the existing <code>Instance</code> .
<code>count_vocab_items(counter, Dict[str, int])</code>	Increments counts in the given <code>counter</code> for all of the vocabulary items in all of the <code>Fields</code> in this <code>Instance</code> .

Continued on next page

Table 8 – continued from previous page

<code>dynamic_index_fields(vocab, dynamic_fields)</code>	<code>dy-</code>	Indexes all fields in this Instance using the provided Vocabulary.
<code>get(k,d)</code>		
<code>index_fields(vocab)</code>		Indexes all fields in this Instance using the provided Vocabulary.
<code>items()</code>		
<code>keys()</code>		
<code>values()</code>		

**add\_field** (*field*: `antu.io.fields.field.Field`) → None  
 Add the field to the existing Instance.

**Parameters**

**field** [`Field`] Which field needs to be added.

**count\_vocab\_items** (*counter*: `Dict[str, Dict[str, int]]`) → None

Increments counts in the given `counter` for all of the vocabulary items in all of the `Fields` in this Instance.

**Parameters**

**counter** [`Dict[str, Dict[str, int]]`] We count the number of strings if the string needs to be counted to some counters.

**dynamic\_index\_fields** (*vocab*: `antu.io.vocabulary.Vocabulary`, *dynamic\_fields*: `Set[str]`) → `Dict[str, Dict[str, Indices]]`

Indexes all fields in this Instance using the provided Vocabulary. This *mutates* the current object, it does not return a new Instance. A `DataIterator` will call this on each pass through a dataset; we use the `indexed` flag to make sure that indexing only happens once. This means that if for some reason you modify your vocabulary after you’ve indexed your instances, you might get unexpected behavior.

**Parameters**

**vocab** [`Vocabulary`] `vocab` is used to get the index of each item.

**Returns**

**res** [`Dict[str, Dict[str, Indices]]`] Returns the Indices corresponding to the instance. The first key is field name and the second key is the vocabulary name.

**index\_fields** (*vocab*: `antu.io.vocabulary.Vocabulary`) → `Dict[str, Dict[str, Indices]]`

Indexes all fields in this Instance using the provided Vocabulary. This *mutates* the current object, it does not return a new Instance. A `DataIterator` will call this on each pass through a dataset; we use the `indexed` flag to make sure that indexing only happens once. This means that if for some reason you modify your vocabulary after you’ve indexed your instances, you might get unexpected behavior.

**Parameters**

**vocab** [`Vocabulary`] `vocab` is used to get the index of each item.

**Returns**

**res** [`Dict[str, Dict[str, Indices]]`] Returns the Indices corresponding to the instance. The first key is field name and the second key is the vocabulary name.

## 1.4 antu.io.vocabulary module

```
class antu.io.vocabulary.Vocabulary (counters: Dict[str, Dict[str, int]] = {}, min_count:
    Dict[str, int] = {}, pretrained_vocab: Dict[str,
    List[str]] = {}, intersection_vocab: Dict[str, str] = {},
    no_pad_namespace: Set[str] = {}, no_unk_namespace:
    Set[str] = {})
```

Bases: `object`

### Parameters

- counters** [Dict[str, Dict[str, int]], optional (default= dict () )] Element statistics for datasets.
- min\_count** [Dict[str, int], optional (default= dict () )] Defines the minimum number of occurrences when some counter are converted to vocabulary.
- pretrained\_vocab** [Dict[str, List[str]], optional (default= dict () )] External pre-trained vocabulary.
- intersection\_vocab** [Dict[str, str], optional (default= dict () )] Defines the intersection with which vocabulary takes, when loading some oversized pre-trained vocabulary.
- no\_pad\_namespace** [Set[str], optional (default= set () )] Defines which vocabularies do not have *pad* token.
- no\_unk\_namespace** [Set[str], optional (default= set () )] Defines which vocabularies do not have *oov* token.

### Methods

<code>add_token_to_namespace(token, namespace)</code>	Extend the vocabulary by add token to vocabulary namespace.
<code>extend_from_counter(counters, Dict[str, ...])</code>	Extend the vocabulary from the dataset statistic counters after defining the vocabulary.
<code>extend_from_pretrained_vocab(...)</code>	Extend the vocabulary from the pre-trained vocabulary after defining the vocabulary.
<code>get_token_from_index(index, vocab_name)</code>	Gets the token of a index in the vocabulary.
<code>get_token_index(token, vocab_name)</code>	Gets the index of a token in the vocabulary.
<code>get_vocab_size(namespace)</code>	Gets the size of a vocabulary.

<code>get_padding_index</code>	
<code>get_unknow_index</code>	

**add\_token\_to\_namespace** (*token: str, namespace: str*) → None  
 Extend the vocabulary by add token to vocabulary namespace.

### Parameters

- token** [str] The token that needs to be added.
  - namespace** [str] Which vocabulary needs to be added to.
- extend\_from\_counter** (*counters: Dict[str, Dict[str, int]], min\_count: Union[int, Dict[str, int]] = {}, no\_pad\_namespace: Set[str] = {}, no\_unk\_namespace: Set[str] = {}*) → None

Extend the vocabulary from the dataset statistic counters after defining the vocabulary.

**Parameters**

**counters** [Dict[str, Dict[str, int]]] Element statistics for datasets.

**min\_count** [Dict[str, int], optional (default= dict() )] Defines the minimum number of occurrences when some counter are converted to vocabulary.

**no\_pad\_namespace** [Set[str], optional (default= set() )] Defines which vocabularies do not have *pad* token.

**no\_unk\_namespace** [Set[str], optional (default= set() )] Defines which vocabularies do not have *oov* token.

**extend\_from\_pretrained\_vocab** (*pretrained\_vocab*: Dict[str, List[str]], *intersection\_vocab*: Dict[str, str] = {}, *no\_pad\_namespace*: Set[str] = {}, *no\_unk\_namespace*: Set[str] = {}) → None

Extend the vocabulary from the pre-trained vocabulary after defining the vocabulary.

**Parameters**

**pretrained\_vocab** [Dict[str, List[str]]] External pre-trained vocabulary.

**intersection\_vocab** [Dict[str, str], optional (default= dict() )] Defines the intersection with which vocabulary takes, when loading some oversized pre-trained vocabulary.

**no\_pad\_namespace** [Set[str], optional (default= set() )] Defines which vocabularies do not have *pad* token.

**no\_unk\_namespace** [Set[str], optional (default= set() )] Defines which vocabularies do not have *oov* token.

**get\_padding\_index** (*namespace*: str) → int

**get\_token\_from\_index** (*index*: int, *vocab\_name*: str) → str

Gets the token of a index in the vocabulary.

**Parameters**

**index** [int] Gets the token of which index.

**namespace** [str] Which vocabulary this index belongs to.

**Returns**

**Token** [str]

**get\_token\_index** (*token*: str, *vocab\_name*: str) → int

Gets the index of a token in the vocabulary.

**Parameters**

**token** [str] Gets the index of which token.

**namespace** [str] Which vocabulary this token belongs to.

**Returns**

**Index** [int]

**get\_unknow\_index** (*namespace*: str) → int

**get\_vocab\_size** (*namespace*: str) → int

Gets the size of a vocabulary.

**Parameters**



**namespace** [str] Which vocabulary.

**Returns**

**Vocabulary size** [int]

## 1.5 Module contents



## 2.1 Subpackages

### 2.1.1 antu.nn.dynet package

#### Submodules

antu.nn.dynet.attention\_mechanism module

antu.nn.dynet.char2word\_embedder module

antu.nn.dynet.initializer module

antu.nn.dynet.multi\_layer\_perception module

antu.nn.dynet.nn\_classifier module

antu.nn.dynet.rnn\_builder module

#### Module contents

## 2.2 Module contents



## CHAPTER 3

---

### Indices and tables

---

- `genindex`
- `modindex`
- `search`



**a**

antu.io, 13  
antu.io.dataset\_readers, 3  
antu.io.dataset\_readers.dataset\_reader,  
    3  
antu.io.datasets, 4  
antu.io.datasets.dataset, 4  
antu.io.fields, 7  
antu.io.fields.field, 4  
antu.io.fields.index\_field, 5  
antu.io.fields.sequence\_label\_field, 5  
antu.io.fields.text\_field, 6  
antu.io.instance, 9  
antu.io.token\_indexers, 9  
antu.io.token\_indexers.char\_token\_indexer,  
    7  
antu.io.token\_indexers.single\_id\_token\_indexer,  
    8  
antu.io.token\_indexers.token\_indexer, 9  
antu.io.vocabulary, 11  
antu.nn, 15





## A

add\_field() (*antu.io.instance.Instance* method), 10  
 add\_token\_to\_namespace() (*antu.io.vocabulary.Vocabulary* method), 11  
 antu.io (*module*), 13  
 antu.io.dataset\_readers (*module*), 3  
 antu.io.dataset\_readers.dataset\_reader (*module*), 3  
 antu.io.datasets (*module*), 4  
 antu.io.datasets.dataset (*module*), 4  
 antu.io.fields (*module*), 7  
 antu.io.fields.field (*module*), 4  
 antu.io.fields.index\_field (*module*), 5  
 antu.io.fields.sequence\_label\_field (*module*), 5  
 antu.io.fields.text\_field (*module*), 6  
 antu.io.instance (*module*), 9  
 antu.io.token\_indexers (*module*), 9  
 antu.io.token\_indexers.char\_token\_indexer (*module*), 7  
 antu.io.token\_indexers.single\_id\_token\_indexer (*module*), 8  
 antu.io.token\_indexers.token\_indexer (*module*), 9  
 antu.io.vocabulary (*module*), 11  
 antu.nn (*module*), 15

## B

build\_dataset() (*antu.io.datasets.dataset.Dataset* method), 4

## C

CharTokenIndexer (*class* in *antu.io.token\_indexers.char\_token\_indexer*), 7  
 count\_vocab\_items() (*antu.io.fields.field.Field* method), 4  
 count\_vocab\_items() (*antu.io.fields.index\_field.IndexField* method),

5

count\_vocab\_items() (*antu.io.fields.sequence\_label\_field.SequenceLabelField* method), 5  
 count\_vocab\_items() (*antu.io.fields.text\_field.TextField* method), 6  
 count\_vocab\_items() (*antu.io.instance.Instance* method), 10  
 count\_vocab\_items() (*antu.io.token\_indexers.char\_token\_indexer.CharTokenIndexer* method), 7  
 count\_vocab\_items() (*antu.io.token\_indexers.single\_id\_token\_indexer.SingleIdTokenIndexer* method), 8  
 count\_vocab\_items() (*antu.io.token\_indexers.token\_indexer.TokenIndexer* method), 9

## D

Dataset (*class* in *antu.io.datasets.dataset*), 4  
 DatasetReader (*class* in *antu.io.dataset\_readers.dataset\_reader*), 3  
 datasets (*antu.io.datasets.dataset.Dataset* attribute), 4  
 dynamic\_index\_fields() (*antu.io.instance.Instance* method), 10

## E

extend\_from\_counter() (*antu.io.vocabulary.Vocabulary* method), 11  
 extend\_from\_pretrained\_vocab() (*antu.io.vocabulary.Vocabulary* method), 12

## F

Field (*class* in *antu.io.fields.field*), 4

## G

get\_padding\_index() (antu.io.vocabulary.Vocabulary method), 12

get\_token\_from\_index() (antu.io.vocabulary.Vocabulary method), 12

get\_token\_index() (antu.io.vocabulary.Vocabulary method), 12

get\_unknow\_index() (antu.io.vocabulary.Vocabulary method), 12

get\_vocab\_size() (antu.io.vocabulary.Vocabulary method), 12

tokens\_to\_indices() (antu.io.token\_indexers.token\_indexer.TokenIndexer method), 9

## V

Vocabulary (class in antu.io.vocabulary), 11

vocabulary\_set (antu.io.datasets.dataset.Dataset attribute), 4

## I

index() (antu.io.fields.field.Field method), 4

index() (antu.io.fields.index\_field.IndexField method), 5

index() (antu.io.fields.sequence\_label\_field.SequenceLabelField method), 6

index() (antu.io.fields.text\_field.TextField method), 6

index\_fields() (antu.io.instance.Instance method), 10

IndexField (class in antu.io.fields.index\_field), 5

input\_to\_instance() (antu.io.dataset\_readers.dataset\_reader.DatasetReader method), 3

Instance (class in antu.io.instance), 9

## R

read() (antu.io.dataset\_readers.dataset\_reader.DatasetReader method), 3

## S

SequenceLabelField (class in antu.io.fields.sequence\_label\_field), 5

SingleIdTokenIndexer (class in antu.io.token\_indexers.single\_id\_token\_indexer), 8

## T

TextField (class in antu.io.fields.text\_field), 6

TokenIndexer (class in antu.io.token\_indexers.token\_indexer), 9

tokens\_to\_indices() (antu.io.token\_indexers.char\_token\_indexer.CharTokenIndexer method), 7

tokens\_to\_indices() (antu.io.token\_indexers.single\_id\_token\_indexer.SingleIdTokenIndexer method), 8