# **IPOP Cabinet**

None

# Table of contents

	Jser Guide	3
	API Reference	4
	1 Getting Started	4
	2 UMLS Drawer	5
2.3	3 Cleaning Drawer	10

## 1. User Guide

\*\*\*

\*\*\*# data-utils

rename to something... cabinet

all sub-folders will be suffixed with "\_drawer"

TODO: switch to mkdocs, boo sphinx

:)

metamap TOS: https://lhncbc.nlm.nih.gov/ii/tools/MetaMap/run-locally/Ts\_and\_Cs.html

does metamap have a "validated/authorized" endpoint?

windows install: https://lhncbc.nlm.nih.gov/ii/tools/MetaMap/Docs/README\_win32.html

Times for metamap vs scispacy...

img

img

These are "straight up" and do not involve parallelization, in that regard:

- metmap benefits about 4x improvement from IO thread-based multiprocessing
- scispacy can utilize batch processing link and can benefit even further from CPU (cpus go zoom) parallelization in with batch processing and limiting pipeline models
- HOWEVER
- a fastapi endpoint (based on internet search) can easily serve about 200req/s which is faster than scispacy model of about 100it/s
- so ok to use endpoint
- we can further unlock this using websockets
- ner endpoint that does the SNOMED traversal
- something like... `post ner(text: str, terminal tree node: str -> cui) ?

future work will turn this into a rust/python package using pyo3 and maturin... this isn't necessary for current system calls and may benefit network calls but will mostly be required for parsing the snomed source files when we support providing your own instead of just (implying we keep current functionality as well) the id-based maps

this will also be supported by rust libraries such as mmi-parser to parse out mmi output from the rest api

## 2. API Reference

## 2.1 Getting Started

This is where you can find source code documentation, examples, and more technical details.

Our cabinet of tools.

#### ghlights 🗸

- Local MetaMap operations
- SciSpacy NER via API
- SNOMED tree traversal
- UMLS CUI to SNOMED CUI
- Common data normalization tasks

In general, we try to expose the high-level functionality of these tools at the top level of their corresponding "drawers" (e.g. cabinet.umls\_drawer).

This way you can import a drawer and use its functionality specifically.

For example:

from cabinet import umls\_drawer umls\_drawer.post\_ner\_single("I have a headache.")

If you want more granular control/exposure, check out the underscore methods inside the drawers although this is not recommended practice.

#### 2.2 UMLS Drawer

This drawer is for UMLS related activities.

The scispacy\_ner module gives you access to the scispacy biomedical NER model via our API.

The metamap\_ner module interacts with the MetaMap NLP tool to extract structured information from biomedical text and requires you to have MetaMap installed locally.

The knowledge\_base module allows you to interact with the UMLS knowledge base data at a high level and mostly focuses on SNOMED CT concepts. Further work on this module may take advantage of the *entire* UMLS, but require a locally downloaded copy due to licensing restrictions.

In general, we recommend using the scispacy\_ner module for NER tasks and the knowledge\_base module for knowledge base related tasks unless you specifically need the power of MetaMap.

The post\_ner methods exposed here utilize the API to perform NER on your text.

This module contains functions for interacting with the scispacy NER model via our API.

The private functions utilize async/await syntax and are used by the public functions which are synchronous. The public functions are the ones that you should use in your code unless you are confident that you know what you are doing.

The core type of this module is NEROutput which is a pydantic model that represents the output from the scispacy NER model. All public functions return either an instance of this class or an iterator of instances of this class attached to an index (tuple[int, NEROutput]) for the index of the text that was submitted... this helps with link to original data.

 ${\tt web\_socket\_ner}\ ,\ specifically,\ returns\ an\ iterator\ and\ thus\ needs\ to\ be\ consumed\ to\ be\ used:$ 

```
>>> from cabinet.umls_drawer.scispacy import web_socket_ner
>>> for text_index, ner_output in web_socket_ner(texts=["cocaine", "heroin", "cociane"]):
... print(text_index, ner_output)
0 cui='12' concept_name='test' concept_definition='test22' entity='{"text":"cocaine"}' score=1.0
1 cui='12' concept_name='test' concept_definition='test22' entity='{"text":"heroin"}' score=1.0
2 cui='12' concept_name='test' concept_definition='test22' entity='{"text":"cociane"}' score=1.0
```

#### 2.2.1 NEROutput

Bases: BaseModel

Output from the (scispacy) NER model.

This class is keyword only so you must pass in the arguments as: cui="C0004096", concept\_name="Acetaminophen", ...

Args/Attributes: cui (str): The UMLS CUI. entity (str): The entity that matched a UMLS concept from the source text. score (float): The score of the match.

#### **Examples:**

An example of manually creating this class:

```
>>> from cabinet.umls_drawer import NEROutput
>>> NEROutput(
... cui="C0004096",
... entity="acetaminophen",
... score=0.96,
...)
```

However, much more likely is that you get this class as a return type from one of the various functions in this module that make calls to our API.

```
>>> from cabinet.umls_drawer import post_ner_single
>>> post_ner_single(text="cocaine")
[
```

```
0 cui='12' entity="cocaine"} score=1.0
]
```

## Source code in src/cabinet/umls\_drawer/scispacy\_ner.py ~ class NEROutput(BaseModel): """Output from the ([scispacy](https://github.com/allenai/scispacy/tree/4f9ba0931d216ddfb9a8f01334d76cfb662738ae)) NER model. 35 36 37 38 39 40 41 42 This class is keyword only so you must pass in the arguments as: `cui="C0004096", concept\_name="Acetaminophen", ...` cui (str): The UMLS CUI. entity (str): The entity that matched a UMLS concept from the source text. score (float): The score of the match. 43 44 45 An example of manually creating this class: ```python >>> from cabinet.umls\_drawer import NEROutput 46 47 48 49 >>> NEROutput( ... cui="C0004096", ... entity="acetaminophen", ... ) score=0.96, 50 51 52 53 54 55 56 57 58 59 However, much more likely is that you get this class as a return type from one of the various functions in this module that ```python >>> from cabinet.umls\_drawer import post\_ner\_single >>> post\_ner\_single(text="cocaine") 0 cui='12' entity="cocaine"} score=1.0 60 61 62 63 64 65 66 cui: str ""The UMLS CUI.""" 68 """The entity that matched a UMLS concept from the source text.""" score: float 70 """The score of the match."""

cui: str class-attribute

The UMLS CUI.

entity: str class-attribute

The entity that matched a UMLS concept from the source text.

score: float class-attribute

The score of the match.

#### 2.2.2 post\_ner\_many(texts)

Submit multiple text blobs to the scispacy NER model and return the results.

#### **Parameters:**

Name	Туре	Description	Default
texts	list[str]	The texts to submit to the NER model.	required

#### **Returns:**

Туре	Description
<pre>list[tuple[int, NEROutput]]</pre>	$Iterator[tuple[int, NEROutput]]: The \ results \ from \ the \ NER \ model.$

#### **Raises:**

Type Description

Exception If the response status is not 200.

## Example Y

#### ```python

from cabinet.umls\_drawer import post\_ner\_many post\_ner\_many(texts=["acetaminophen", "ibuprofen"]) [ (0, NEROutput(cui="C0004096", entity="acetaminophen", score=0.96, )), (1, NEROutput(cui="C0004096", entity="ibuprofen", score=0.96, )) ]

#### Source code in src/cabinet/umls\_drawer/scispacy\_ner.py ~ 174 @validate\_arguments def post\_ner\_many( texts: list[str], 175 176 177 ) -> list[tuple[int, NEROutput]]: """Submit multiple text blobs to the scispacy NER model and return the results. 178 179 180 texts (list[str]): The texts to submit to the NER model. 181 182 183 184 Iterator[tuple[int, NEROutput]]: The results from the NER model. 185 186 187 188 189 Exception: If the response status is not 200. "`python >>> from cabinet.umls\_drawer import post\_ner\_many >>> post\_ner\_many(texts=["acetaminophen", "ibuprofen"]) 190 191 192 193 194 (0, NEROutput( cui="C0004096", entity="acetaminophen", score=0.96, 195 196 197 198 199 )), (1, NEROutput( 200 201 cui="C0004096" entity="ibuprofen", score=0.96, 202 204 205 return asyncio.run(\_post\_ner\_many(texts)) 206

#### 2.2.3 post\_ner\_single(text)

Submit a single text blob to the scispacy NER model and return the results.

#### **Parameters:**

Name	Туре	Description	Default
text	str	The text to submit to the NER model.	required

#### **Returns:**

Туре	Description
list[NEROutput]	$list[{\tt NEROutput}]{\tt :}\ {\tt The\ results\ from\ the\ NER\ model}.$

#### Raises:

Туре	Description
Exception	If the response status is not 200.

## Example Y

## ```python

from cabinet.umls\_drawer import post\_ner\_single post\_ner\_single(text="acetaminophen") [ NEROutput( cui="C0004096", entity="acetaminophen", score=0.96, ) ]

#### Source code in src/cabinet/umls\_drawer/scispacy\_ner.py ~ @validate\_arguments def post\_ner\_single(text: str) -> list[NEROutput]: """Submit a single text blob to the scispacy NER model and return the results. 147 148 149 150 Args: text (str): The text to submit to the NER model. 151 152 153 list[NEROutput]: The results from the NER model. 155 156 157 Exception: If the response status is not 200. 158 159 Example: '``python >>> from cabinet.umls\_drawer import post\_ner\_single >>> post\_ner\_single(text="acetaminophen") 160 161 162 163 NEROutput( 164 165 cui="C0004096", entity="acetaminophen", 166 167 score=0.96, 168 170 171 ${\tt return \ asyncio.run(\_post\_nlp\_single(text))}$

#### 2.2.4 websocket\_ner(texts) async

Connect to the scispacy NER model websocket and submit texts.

*IMPORTANT*: This function requires using the async for syntax and thus may not work in all scenarios or environments. It exists for **very** large datasets where the overhead of the HTTP request/response cycle is too much.

#### **Parameters:**

Name	Туре	Description	Default
texts	list[str]	The texts to submit to the NER model.	required

#### Yields:

Туре	Description
AsyncIterator[tuple[int, NEROutput]]	AsyncIterator[tuple[int, NEROutput]]: The results from the NER model.

#### Raises:

Туре	Description
Exception	If the response status is not 200.

## Example Y

#### ```python

from cabinet.umls\_drawer import websocket\_ner async for i, result in websocket\_ner(texts=["acetaminophen", "ibuprofen"]): ... print(i, result) 0 cui='12' concept\_name='test' entity="cocaine" score=1.0 1 cui='12' concept\_name='test' entity="heroin" score=1.0 2 cui='12' concept\_name='test' entity="cociane" score=1.0

## Source code in src/cabinet/umls\_drawer/scispacy\_ner.py

```
async def websocket_ner(texts: list[str]) -> AsyncIterator[tuple[int, NEROutput]]:
"""Connect to the scispacy NER model websocket and submit texts.
213
214
                *IMPORTANT*: This function requires using the `async for` syntax and thus may not work in all scenarios or environments. It exists for **very** large datasets where the overhead of the HTTP request/response cycle is too much.
215
217
218
                         texts (list[str]): The texts to submit to the NER model.
219
220
221
                       AsyncIterator[tuple[int, NEROutput]]: The results from the NER model.
223
224
225
                        Exception: If the response status is not 200.
                Example:
    ```python
    >>> from cabinet.umls_drawer import websocket_ner
226
227
228
229
                        >>> async for i, result in websocket_ner(texts=["acetaminophen", "ibuprofen"]):
... print(i, result)
230
231
                      o cui='12' concept_name='test' entity="cocaine" score=1.0

cui='12' concept_name='test' entity="heroin" score=1.0

cui='12' concept_name='test' entity="cociane" score=1.0
232
233
234
                 async with aiohttp.ClientSession(_WS_URL) as session:
    async with session.ws_connect("/models/ner/ws") as ws:
236
237
                               for i, text in tqdm_asyncio(enumerate(texts)):
    await ws.send_bytes(orjson.dumps({"text": text}))
    raw_data = await ws.receive_bytes()
    response_data = orjson.loads(raw_data)
    data = NEROutput(**response_data)
238
240
                                       yield (i, data)
242
```

## 2.3 Cleaning Drawer

This drawer exists for cleaning functions.

It's a bit of a catch-all for functions that don't fit in any other drawer.

This module contains code for typical data normalization tasks.

 $We generally \ enforce \ 'type-saftey' \ by \ using \ pydantic's \ validate\_arguments \ decorator \ and \ other \ pydantic \ types.$ 

## 2.3.1 categorize\_age(age)

Categorize an age into a string.

#### **Parameters:**

Name	Туре	Description	Default
age	PositiveInt	Age to categorize. Must be a positive integer.	required

#### **Returns:**

Name	Type	Description
str	str	Categorized age.

## **Examples:**

```
>>> categorize_age(10)
'<18'
>>> categorize_age(18)
'18-25'
>>> # a general example using pandas
>>> df['age'].apply(categorize_age)
pandas.Series(['<18', '18-25', '26-35', '36-45', '46-55', '56-65', '65+'])
```

```
## Source code in src/cabinet/cleaning_drawer/normalize.py

## ## Source code in src/cabinet/cleaning_drawer/cabinet/cleaning_drawer/cabinet
```