```
1 !pip install langchain==0.0.267
2 !pip install requests
3 !pip install BeautifulSoup
```

Requirement already satisfied: greenlet>=1 in /usr/local/lib/python3.11/dist-packages (from SQLAlchemy<3,>=1.4->langchain==0.0.267) (3.1.1)
Requirement already satisfied: packaging>=17.0 in /usr/local/lib/python3.11/dist-packages (from marshmallow<4.0.0,>=3.18.0->dataclasses-json<0.6.0,>=0.5.7->langchain
Collecting mypy-extensions>=0.3.0 (from typing-inspect<1,>=0.4.0->dataclasses-json<0.6.0,>=0.5.7->langchain==0.0.267)
  Downloading mypy_extensions-1.0.0-py3-none-any.whl.metadata (1.1 kB)
Downloading langchain-0.0.267-py3-none-any.whl (1.5 MB)
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 1.5/1.5 MB 35.1 MB/s eta 0:00:00
Downloading dataclasses_json-0.5.14-py3-none-any.whl (26 kB)
Downloading langsmith-0.0.92-py3-none-any.whl (56 kB)
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 56.5/56.5 kB 4.5 MB/s eta 0:00:00

```
  × python setup.py egg_info did not run successfully.
  │ exit code: 1
  ╰─> See above for output.

  note: This error originates from a subprocess, and is likely not a problem with pip.
  Preparing metadata (setup.py) ... error
error: metadata-generation-failed

  × Encountered error while generating package metadata.
```

```
1 !pip uninstall spacy cymem murmurhash preshed thinc blis -y
2 !pip install spacy==3.5.0
3 !pip install -numpy
4 #!spacy.prefer_gpu()
5 !pip install scispacy
```

```
Found existing installation: spacy 3.8.5
Uninstalling spacy-3.8.5:
  Successfully uninstalled spacy-3.8.5
Found existing installation: cymem 2.0.11
Uninstalling cymem-2.0.11:
  Successfully uninstalled cymem-2.0.11
Found existing installation: murmurhash 1.0.12
Uninstalling murmurhash-1.0.12:
  Successfully uninstalled murmurhash-1.0.12
Found existing installation: preshed 3.0.9
Uninstalling preshed-3.0.9:
  Successfully uninstalled preshed-3.0.9
Found existing installation: thinc 8.3.6
Uninstalling thinc-8.3.6:
  Successfully uninstalled thinc-8.3.6
Found existing installation: blis 1.3.0
Uninstalling blis-1.3.0:
  Successfully uninstalled blis-1.3.0
Collecting spacy==3.5.0
  Downloading spacy-3.5.0-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (25 kB)
Requirement already satisfied: spacy-legacy<3.1.0,>=3.0.11 in /usr/local/lib/python3.11/dist-packages (from spacy==3.5.0) (3.0.12)
Requirement already satisfied: spacy-loggers<2.0.0,>=1.0.0 in /usr/local/lib/python3.11/dist-packages (from spacy==3.5.0) (1.0.5)
Collecting murmurhash<1.1.0,>=0.28.0 (from spacy==3.5.0)
  Downloading murmurhash-1.0.12-cp311-cp311-manylinux_2_5_x86_64.manylinux1_x86_64.manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (2.1 kB)
Collecting cymem<2.1.0,>=2.0.2 (from spacy==3.5.0)
  Downloading cymem-2.0.11-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (8.5 kB)
Collecting preshed<3.1.0,>=3.0.2 (from spacy==3.5.0)
  Downloading preshed-3.0.9-cp311-cp311-manylinux_2_5_x86_64.manylinux1_x86_64.manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (2.2 kB)
Collecting thinc<8.2.0,>=8.1.0 (from spacy==3.5.0)
  Downloading thinc-8.1.12-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (15 kB)
Requirement already satisfied: wasabi<1.2.0,>=0.9.1 in /usr/local/lib/python3.11/dist-packages (from spacy==3.5.0) (1.1.3)
Requirement already satisfied: srsly<3.0.0,>=2.4.3 in /usr/local/lib/python3.11/dist-packages (from spacy==3.5.0) (2.5.1)
Requirement already satisfied: catalogue<2.1.0,>=2.0.6 in /usr/local/lib/python3.11/dist-packages (from spacy==3.5.0) (2.0.10)
Collecting typer<0.8.0,>=0.3.0 (from spacy==3.5.0)
  Downloading typer-0.7.0-py3-none-any.whl.metadata (17 kB)
Collecting pathy>=0.10.0 (from spacy==3.5.0)
  Downloading pathy-0.11.0-py3-none-any.whl.metadata (16 kB)
Collecting smart-open<7.0.0,>=5.2.1 (from spacy==3.5.0)
  Downloading smart_open-6.4.0-py3-none-any.whl.metadata (21 kB)
Requirement already satisfied: tqdm<5.0.0,>=4.38.0 in /usr/local/lib/python3.11/dist-packages (from spacy==3.5.0) (4.67.1)
Requirement already satisfied: numpy>=1.15.0 in /usr/local/lib/python3.11/dist-packages (from spacy==3.5.0) (1.26.4)
Requirement already satisfied: requests<3.0.0,>=2.13.0 in /usr/local/lib/python3.11/dist-packages (from spacy==3.5.0) (2.32.3)
Collecting pydantic!=1.8,!=1.8.1,<1.11.0,>=1.7.4 (from spacy==3.5.0)
  Downloading pydantic-1.10.21-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (153 kB)
  ━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 153.9/153.9 kB 10.2 MB/s eta 0:00:00
Requirement already satisfied: jinja2 in /usr/local/lib/python3.11/dist-packages (from spacy==3.5.0) (3.1.6)
Requirement already satisfied: setuptools in /usr/local/lib/python3.11/dist-packages (from spacy==3.5.0) (75.2.0)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from spacy==3.5.0) (24.2)
Requirement already satisfied: langcodes<4.0.0,>=3.2.0 in /usr/local/lib/python3.11/dist-packages (from spacy==3.5.0) (3.5.0)
Requirement already satisfied: language-data>=1.2 in /usr/local/lib/python3.11/dist-packages (from langcodes<4.0.0,>=3.2.0->spacy==3.5.0) (1.3.0)
Collecting pathlib-abc==0.1.1 (from pathy>=0.10.0->spacy==3.5.0)
  Downloading pathlib_abc-0.1.1-py3-none-any.whl.metadata (18 kB)
Requirement already satisfied: typing-extensions>=4.2.0 in /usr/local/lib/python3.11/dist-packages (from pydantic!=1.8,!=1.8.1,<1.11.0,>=1.7.4->spacy==3.5.0) (4.13.1
```

```
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests<3.0.0,>=2.13.0->spacy==3.5.0) (3.4.1)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests<3.0.0,>=2.13.0->spacy==3.5.0) (3.10)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests<3.0.0,>=2.13.0->spacy==3.5.0) (2.3.0)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests<3.0.0,>=2.13.0->spacy==3.5.0) (2025.1.31)
Collecting blis<0.8.0,>=0.7.8 (from thinc<8.2.0,>=8.1.0->spacy==3.5.0)
  Downloading blis-0.7.11-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (7.4 kB)
Requirement already satisfied: confection<1.0.0,>=0.0.1 in /usr/local/lib/python3.11/dist-packages (from thinc<8.2.0,>=8.1.0->spacy==3.5.0) (0.1.5)
Requirement already satisfied: click<9.0.0,>=7.1.1 in /usr/local/lib/python3.11/dist-packages (from typer<0.8.0,>=0.3.0->spacy==3.5.0) (8.1.8)
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.11/dist-packages (from jinja2->spacy==3.5.0) (3.0.2)
Requirement already satisfied: marisa-trie>=1.1.0 in /usr/local/lib/python3.11/dist-packages (from language-data>=1.2->langcodes<4.0.0,>=3.2.0->spacy==3.5.0) (1.2.1)
Downloading spacy-3.5.0-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (6.6 MB)
   ──────────────────────────────────────── 6.6/6.6 MB 12.0 MB/s eta 0:00:00
Downloading cymem-2.0.11-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (218 kB)
   ──────────────────────────────────────── 218.9/218.9 kB 12.7 MB/s eta 0:00:00
Downloading murmurhash-1.0.12-cp311-cp311-manylinux_2_5_x86_64.manylinux1_x86_64.manylinux_2_17_x86_64.manylinux2014_x86_64.whl (134 kB)
   ──────────────────────────────────────── 134.3/134.3 kB 9.7 MB/s eta 0:00:00
Downloading pathy-0.11.0-py3-none-any.whl (47 kB)
   ──────────────────────────────────────── 47.3/47.3 kB 3.2 MB/s eta 0:00:00
Downloading pathlib_abc-0.1.1-py3-none-any.whl (23 kB)
Downloading preshed-3.0.9-cp311-cp311-manylinux_2_5_x86_64.manylinux1_x86_64.manylinux_2_17_x86_64.manylinux2014_x86_64.whl (157 kB)
   ──────────────────────────────────────── 157.2/157.2 kB 8.5 MB/s eta 0:00:00
Downloading pydantic-1.10.21-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (3.1 MB)
   ──────────────────────────────────────── 3.1/3.1 MB 63.7 MB/s eta 0:00:00
Downloading smart_open-6.4.0-py3-none-any.whl (57 kB)
   ──────────────────────────────────────── 57.0/57.0 kB 4.1 MB/s eta 0:00:00
Downloading thinc-8.1.12-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (917 kB)
   ──────────────────────────────────────── 917.4/917.4 kB 38.8 MB/s eta 0:00:00
Downloading typer-0.7.0-py3-none-any.whl (38 kB)
Downloading blis-0.7.11-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (10.2 MB)
   ──────────────────────────────────────── 10.2/10.2 MB 63.3 MB/s eta 0:00:00
Installing collected packages: cymem, typer, smart-open, pydantic, pathlib-abc, murmurhash, blis, preshed, pathy, thinc, spacy
  Attempting uninstall: typer
    Found existing installation: typer 0.15.2
    Uninstalling typer-0.15.2:
      Successfully uninstalled typer-0.15.2
  Attempting uninstall: smart-open
    Found existing installation: smart-open 7.1.0
    Uninstalling smart-open-7.1.0:
      Successfully uninstalled smart-open-7.1.0
  Attempting uninstall: pydantic
    Found existing installation: pydantic 2.11.2
    Uninstalling pydantic-2.11.2:
      Successfully uninstalled pydantic-2.11.2
ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour is the source of the following dependency c
langchain-core 0.3.51 requires langsmith<0.4,>=0.1.125, but you have langsmith 0.0.92 which is incompatible.
langchain-core 0.3.51 requires pydantic<3.0.0,>=2.5.2; python_full_version < "3.12.4", but you have pydantic 1.10.21 which is incompatible.
albumentations 2.0.5 requires pydantic>=2.9.2, but you have pydantic 1.10.21 which is incompatible.
google-genai 1.9.0 requires pydantic<3.0.0,>=2.0.0, but you have pydantic 1.10.21 which is incompatible.
Successfully installed blis-0.7.11 cymem-2.0.11 murmurhash-1.0.12 pathlib-abc-0.1.1 pathy-0.11.0 preshed-3.0.9 pydantic-1.10.21 smart-open-6.4.0 spacy-3.5.0 thinc-8.
WARNING: The following packages were previously imported in this runtime:
  [blis,cymem,murmurhash,preshed,pydantic,spacy,thinc]
You must restart the runtime in order to use newly installed versions.

    RESTART SESSION
```

```
RESTART SESSION

Usage:
  pip3 install [options] <requirement specifier> [package-index-options] ...
  pip3 install [options] -r <requirements file> [package-index-options] ...
  pip3 install [options] [-e] <vcs project url> ...
  pip3 install [options] [-e] <local project path> ...
  pip3 install [options] <archive url/path> ...

no such option: -n
Collecting scispacy
  Downloading scispacy-0.5.5-py3-none-any.whl.metadata (18 kB)
Collecting spacy<3.8.0,>=3.7.0 (from scispacy)
  Downloading spacy-3.7.5-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (27 kB)
Requirement already satisfied: scipy in /usr/local/lib/python3.11/dist-packages (from scispacy) (1.14.1)
Requirement already satisfied: requests<3.0.0,>=2.0.0 in /usr/local/lib/python3.11/dist-packages (from scispacy) (2.32.3)
Collecting conllu (from scispacy)
  Downloading conllu-6.0.0-py3-none-any.whl.metadata (21 kB)
Requirement already satisfied: numpy in /usr/local/lib/python3.11/dist-packages (from scispacy) (1.26.4)
Requirement already satisfied: joblib in /usr/local/lib/python3.11/dist-packages (from scispacy) (1.4.2)
Requirement already satisfied: scikit-learn>=0.20.3 in /usr/local/lib/python3.11/dist-packages (from scispacy) (1.6.1)
Collecting pysbd (from scispacy)
  Downloading pysbd-0.3.4-py3-none-any.whl.metadata (6.1 kB)
Collecting nmslib-metabrainz==2.1.3 (from scispacy)
  Downloading nmslib_metabrainz-2.1.3-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (956 bytes)
Collecting pybind11>=2.2.3 (from nmslib-metabrainz==2.1.3->scispacy)
  Downloading pybind11-2.13.6-py3-none-any.whl.metadata (9.5 kB)
Requirement already satisfied: psutil in /usr/local/lib/python3.11/dist-packages (from nmslib-metabrainz==2.1.3->scispacy) (5.9.5)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests<3.0.0,>=2.0.0->scispacy) (3.4.1)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests<3.0.0,>=2.0.0->scispacy) (3.10)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests<3.0.0,>=2.0.0->scispacy) (2.3.0)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests<3.0.0,>=2.0.0->scispacy) (2025.1.31)
Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn>=0.20.3->scispacy) (3.6.0)
Requirement already satisfied: spacy-legacy<3.1.0,>=3.0.11 in /usr/local/lib/python3.11/dist-packages (from spacy<3.8.0,>=3.7.0->scispacy) (3.0.12)
Requirement already satisfied: spacy-loggers<2.0.0,>=1.0.0 in /usr/local/lib/python3.11/dist-packages (from spacy<3.8.0,>=3.7.0->scispacy) (1.0.5)
Requirement already satisfied: murmurhash<1.1.0,>=0.28.0 in /usr/local/lib/python3.11/dist-packages (from spacy<3.8.0,>=3.7.0->scispacy) (1.0.12)
Requirement already satisfied: cymem<2.1.0,>=2.0.2 in /usr/local/lib/python3.11/dist-packages (from spacy<3.8.0,>=3.7.0->scispacy) (2.0.11)
Requirement already satisfied: preshed<3.1.0,>=3.0.2 in /usr/local/lib/python3.11/dist-packages (from spacy<3.8.0,>=3.7.0->scispacy) (3.0.9)
Collecting thinc<8.3.0,>=8.2.2 (from spacy<3.8.0,>=3.7.0->scispacy)
  Downloading thinc-8.2.5-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (15 kB)
Requirement already satisfied: wasabi<1.2.0,>=0.9.1 in /usr/local/lib/python3.11/dist-packages (from spacy<3.8.0,>=3.7.0->scispacy) (1.1.3)
Requirement already satisfied: srsly<3.0.0,>=2.4.3 in /usr/local/lib/python3.11/dist-packages (from spacy<3.8.0,>=3.7.0->scispacy) (2.5.1)
Requirement already satisfied: catalogue<2.1.0,>=2.0.6 in /usr/local/lib/python3.11/dist-packages (from spacy<3.8.0,>=3.7.0->scispacy) (2.0.10)
Requirement already satisfied: weasel<0.5.0,>=0.1.0 in /usr/local/lib/python3.11/dist-packages (from spacy<3.8.0,>=3.7.0->scispacy) (0.4.1)
Requirement already satisfied: typer<1.0.0,>=0.3.0 in /usr/local/lib/python3.11/dist-packages (from spacy<3.8.0,>=3.7.0->scispacy) (0.7.0)
Requirement already satisfied: tqdm<5.0.0,>=4.38.0 in /usr/local/lib/python3.11/dist-packages (from spacy<3.8.0,>=3.7.0->scispacy) (4.67.1)
Requirement already satisfied: pydantic!=1.8,!=1.8.1,<3.0.0,>=1.7.4 in /usr/local/lib/python3.11/dist-packages (from spacy<3.8.0,>=3.7.0->scispacy) (1.10.21)
Requirement already satisfied: jinja2 in /usr/local/lib/python3.11/dist-packages (from spacy<3.8.0,>=3.7.0->scispacy) (3.1.6)
Requirement already satisfied: setuptools in /usr/local/lib/python3.11/dist-packages (from spacy<3.8.0,>=3.7.0->scispacy) (75.2.0)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from spacy<3.8.0,>=3.7.0->scispacy) (24.2)
Requirement already satisfied: langcodes<4.0.0,>=3.2.0 in /usr/local/lib/python3.11/dist-packages (from spacy<3.8.0,>=3.7.0->scispacy) (3.5.0)
Requirement already satisfied: language-data>=1.2 in /usr/local/lib/python3.11/dist-packages (from langcodes<4.0.0,>=3.2.0->spacy<3.8.0,>=3.7.0->scispacy) (1.3.0)
Requirement already satisfied: typing-extensions>=4.2.0 in /usr/local/lib/python3.11/dist-packages (from pydantic!=1.8,!=1.8.1,<3.0.0,>=1.7.4->spacy<3.8.0,>=3.7.0->s
```

Requirement already satisfied: blis<0.8.0,>=0.7.8 in /usr/local/lib/python3.11/dist-packages (from thinc<8.3.0,>=8.2.2->spacy<3.8.0,>=3.7.0->scispacy) (0.7.11)
Requirement already satisfied: confection<1.0.0,>=0.0.1 in /usr/local/lib/python3.11/dist-packages (from thinc<8.3.0,>=8.2.2->spacy<3.8.0,>=3.7.0->scispacy) (0.1.5)
Requirement already satisfied: click<9.0.0,>=7.1.1 in /usr/local/lib/python3.11/dist-packages (from typer<1.0.0,>=0.3.0->spacy<3.8.0,>=3.7.0->scispacy) (8.1.8)
Requirement already satisfied: cloudpathlib<1.0.0,>=0.7.0 in /usr/local/lib/python3.11/dist-packages (from weasel<0.5.0,>=0.1.0->spacy<3.8.0,>=3.7.0->scispacy) (0.21
Requirement already satisfied: smart-open<8.0.0,>=5.2.1 in /usr/local/lib/python3.11/dist-packages (from weasel<0.5.0,>=0.1.0->spacy<3.8.0,>=3.7.0->scispacy) (6.4.0)
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.11/dist-packages (from jinja2->spacy<3.8.0,>=3.7.0->scispacy) (3.0.2)
Requirement already satisfied: marisa-trie>=1.1.0 in /usr/local/lib/python3.11/dist-packages (from language-data>=1.2->langcodes<4.0.0,>=3.2.0->spacy<3.8.0,>=3.7.0->
Downloading scispacy-0.5.5-py3-none-any.whl (46 kB)
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 46.2/46.2 kB 3.2 MB/s eta 0:00:00
Downloading nmslib_metabrainz-2.1.3-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (14.1 MB)
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 14.1/14.1 MB 58.4 MB/s eta 0:00:00
Downloading spacy-3.7.5-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (6.6 MB)
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 6.6/6.6 MB 64.4 MB/s eta 0:00:00
Downloading conllu-6.0.0-py3-none-any.whl (16 kB)
Downloading pysbd-0.3.4-py3-none-any.whl (71 kB)
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 71.1/71.1 kB 5.0 MB/s eta 0:00:00
Downloading pybind11-2.13.6-py3-none-any.whl (243 kB)
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 243.3/243.3 kB 16.7 MB/s eta 0:00:00
Downloading thinc-8.2.5-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (920 kB)
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 920.2/920.2 kB 36.9 MB/s eta 0:00:00
Installing collected packages: pysbd, pybind11, conllu, nmslib-metabrainz, thinc, spacy, scispacy
  Attempting uninstall: thinc
    Found existing installation: thinc 8.1.12
    Uninstalling thinc-8.1.12:
      Successfully uninstalled thinc-8.1.12
  Attempting uninstall: spacy
    Found existing installation: spacy 3.5.0
    Uninstalling spacy-3.5.0:
      Successfully uninstalled spacy-3.5.0
Successfully installed conllu-6.0.0 nmslib-metabrainz-2.1.3 pybind11-2.13.6 pysbd-0.3.4 scispacy-0.5.5 spacy-3.7.5 thinc-8.2.5
**WARNING: The following packages were previously imported in this runtime:**
  **[spacy,thinc]**
**You must restart the runtime in order to use newly installed versions.**

RESTART SESSION

```
1 !pip install https://s3-us-west-2.amazonaws.com/ai2-s2-scispacy/releases/v0.5.1/en_core_sci_md-0.5.1.tar.gz
```

```
Collecting https://s3-us-west-2.amazonaws.com/ai2-s2-scispacy/releases/v0.5.1/en_core_sci_md-0.5.1.tar.gz
  Downloading https://s3-us-west-2.amazonaws.com/ai2-s2-scispacy/releases/v0.5.1/en_core_sci_md-0.5.1.tar.gz (120.2 MB)
                                    ━━━━━━━━━━━━━━━━━━━━━━━━━ 120.2/120.2 MB 6.5 MB/s eta 0:00:00
  Preparing metadata (setup.py) ... done
Collecting spacy<3.5.0,>=3.4.1 (from en_core_sci_md==0.5.1)
  Downloading spacy-3.4.4-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (24 kB)
Requirement already satisfied: spacy-legacy<3.1.0,>=3.0.10 in /usr/local/lib/python3.11/dist-packages (from spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (3.0.12)
Requirement already satisfied: spacy-loggers<2.0.0,>=1.0.0 in /usr/local/lib/python3.11/dist-packages (from spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (1.0.5)
Requirement already satisfied: murmurhash<1.1.0,>=0.28.0 in /usr/local/lib/python3.11/dist-packages (from spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (1.0.12)
Requirement already satisfied: cymem<2.1.0,>=2.0.2 in /usr/local/lib/python3.11/dist-packages (from spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (2.0.11)
Requirement already satisfied: preshed<3.1.0,>=3.0.2 in /usr/local/lib/python3.11/dist-packages (from spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (3.0.9)
Collecting thinc<8.2.0,>=8.1.0 (from spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1)
  Using cached thinc-8.1.12-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (15 kB)
Collecting wasabi<1.1.0,>=0.9.1 (from spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1)
  Downloading wasabi-0.10.1-py3-none-any.whl.metadata (28 kB)
Requirement already satisfied: srsly<3.0.0,>=2.4.3 in /usr/local/lib/python3.11/dist-packages (from spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (2.5.1)
Requirement already satisfied: catalogue<2.1.0,>=2.0.6 in /usr/local/lib/python3.11/dist-packages (from spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (2.0.10)
Requirement already satisfied: typer<0.8.0,>=0.3.0 in /usr/local/lib/python3.11/dist-packages (from spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (0.7.0)
Requirement already satisfied: pathy>=0.3.5 in /usr/local/lib/python3.11/dist-packages (from spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (0.11.0)
Requirement already satisfied: smart-open<7.0.0,>=5.2.1 in /usr/local/lib/python3.11/dist-packages (from spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (6.4.0)
Requirement already satisfied: tqdm<5.0.0,>=4.38.0 in /usr/local/lib/python3.11/dist-packages (from spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (4.67.1)
Requirement already satisfied: numpy>=1.15.0 in /usr/local/lib/python3.11/dist-packages (from spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (1.26.4)
Requirement already satisfied: requests<3.0.0,>=2.13.0 in /usr/local/lib/python3.11/dist-packages (from spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (2.32.3)
Requirement already satisfied: pydantic!=1.8,!=1.8.1,<1.11.0,>=1.7.4 in /usr/local/lib/python3.11/dist-packages (from spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (1.
Requirement already satisfied: jinja2 in /usr/local/lib/python3.11/dist-packages (from spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (3.1.6)
Requirement already satisfied: setuptools in /usr/local/lib/python3.11/dist-packages (from spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (75.2.0)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (24.2)
Requirement already satisfied: langcodes<4.0.0,>=3.2.0 in /usr/local/lib/python3.11/dist-packages (from spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (3.5.0)
Requirement already satisfied: language-data>=1.2 in /usr/local/lib/python3.11/dist-packages (from langcodes<4.0.0,>=3.2.0->spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.
Requirement already satisfied: pathlib-abc==0.1.1 in /usr/local/lib/python3.11/dist-packages (from pathy>=0.3.5->spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (0.1.1)
Requirement already satisfied: typing-extensions>=4.2.0 in /usr/local/lib/python3.11/dist-packages (from pydantic!=1.8,!=1.8.1,<1.11.0,>=1.7.4->spacy<3.5.0,>=3.4.1->
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests<3.0.0,>=2.13.0->spacy<3.5.0,>=3.4.1->en_core_sci_md
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests<3.0.0,>=2.13.0->spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (3.
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests<3.0.0,>=2.13.0->spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests<3.0.0,>=2.13.0->spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.
Requirement already satisfied: blis<0.8.0,>=0.7.8 in /usr/local/lib/python3.11/dist-packages (from thinc<8.2.0,>=8.1.0->spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (
Requirement already satisfied: confection<1.0.0,>=0.0.1 in /usr/local/lib/python3.11/dist-packages (from thinc<8.2.0,>=8.1.0->spacy<3.5.0,>=3.4.1->en_core_sci_md==0.
Requirement already satisfied: click<9.0.0,>=7.1.1 in /usr/local/lib/python3.11/dist-packages (from typer<0.8.0,>=0.3.0->spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1)
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.11/dist-packages (from jinja2->spacy<3.5.0,>=3.4.1->en_core_sci_md==0.5.1) (3.0.2)
Requirement already satisfied: marisa-trie>=1.1.0 in /usr/local/lib/python3.11/dist-packages (from language-data>=1.2->langcodes<4.0.0,>=3.2.0->spacy<3.5.0,>=3.4.1->
Downloading spacy-3.4.4-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (6.4 MB)
                                    ━━━━━━━━━━━━━━━━━━━━━━━━━ 6.4/6.4 MB 31.3 MB/s eta 0:00:00
Using cached thinc-8.1.12-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (917 kB)
Downloading wasabi-0.10.1-py3-none-any.whl (26 kB)
Building wheels for collected packages: en_core_sci_md
  Building wheel for en_core_sci_md (setup.py) ... done
  Created wheel for en_core_sci_md: filename=en_core_sci_md-0.5.1-py3-none-any.whl size=120253138 sha256=792bd41c7595fd08056d199a4eb1fd65f63750e39961e96c98e4bcb88d53
  Stored in directory: /root/.cache/pip/wheels/0a/50/82/7547d452aa8d5a653fb1271c38113de20f7842effc4b7313d0
Successfully built en_core_sci_md
Installing collected packages: wasabi, thinc, spacy, en_core_sci_md
  Attempting uninstall: wasabi
    Found existing installation: wasabi 1.1.3
    Uninstalling wasabi-1.1.3:
```

```
       Successfully uninstalled wasabi-1.1.3
    Attempting uninstall: thinc
       Found existing installation: thinc 8.2.5
       Uninstalling thinc-8.2.5:
          Successfully uninstalled thinc-8.2.5
    Attempting uninstall: spacy
       Found existing installation: spacy 3.7.5
       Uninstalling spacy-3.7.5:
          Successfully uninstalled spacy-3.7.5
ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour is the source of the following dependency c
scispacy 0.5.5 requires spacy<3.8.0,>=3.7.0, but you have spacy 3.4.4 which is incompatible.
Successfully installed en_core_sci_md-0.5.1 spacy-3.4.4 thinc-8.1.12 wasabi-0.10.1
WARNING: The following packages were previously imported in this runtime:
   [spacy,thinc,wasabi]
You must restart the runtime in order to use newly installed versions.
```

RESTART SESSION

```
1 import os
2 import requests
3 import xml.etree.ElementTree as ET
4 import spacy
5 import torch
6 import torch.nn as nn
7 from transformers import AutoModel, AutoTokenizer
```

```
1 import requests
2 from bs4 import BeautifulSoup
3
4 def download_arxiv_html(arxiv_html_url, save_path):
5     response = requests.get(arxiv_html_url)
6     response.raise_for_status()
7
8     soup = BeautifulSoup(response.text, 'html.parser')
9     for script in soup(['script', 'style']):
10        script.extract()
11
12    cleaned_html = soup.prettify()
13
14    with open(save_path, 'w', encoding='utf-8') as file:
15        file.write(cleaned_html)
16
17 def extract_abstract(html_content):
18    soup = BeautifulSoup(html_content, 'html.parser')
19    abstract_div = soup.find('div', class_='ltx_abstract')
20    abstract = abstract_div.get_text(strip=True) if abstract_div else ""
21    if abstract_div:
22        abstract_div.extract()  # remove abstract from soup
23    return abstract[8:], soup.prettify()
24
25
26
27
28 from bs4 import BeautifulSoup
29
30 # def extract_main_paper_from_html(html_content):
31 #     soup = BeautifulSoup(html_content, 'html.parser')
32 #     main_section = soup.find(id='S1')
33
34 #     if not main_section:
35 #         return ""
36
37 #     # Capture everything starting from the main section
38 #     main_paper_parts = []
```

```
39 #      current = main_section
40 #      while current:
41 #          main_paper_parts.append(str(current))
42 #          current = current.find_next_sibling()
43
44 #      return "\n".join(main_paper_parts)
45
46 from bs4 import BeautifulSoup
47
48 def extract_main_paper_text_from_html(html_content):
49     soup = BeautifulSoup(html_content, 'html.parser')
50     main_section = soup.find(id='S1')
51
52     if not main_section:
53         return ""
54
55     # Extract plain text from main section and its siblings
56     main_text_parts = []
57     current = main_section
58     while current:
59         main_text_parts.append(current.get_text(separator=" ", strip=True))
60         current = current.find_next_sibling()
61
62     return "\n".join(main_text_parts)
63
64
65
66 import re
67 from nltk import WordNetLemmatizer
68 import html
69
70 import re
71 import html
72 from bs4 import BeautifulSoup
73 import nltk
74 nltk.download('stopwords')
75 from nltk.tokenize import sent_tokenize, word_tokenize
76 from nltk.corpus import stopwords
77 from nltk.stem import WordNetLemmatizer
78 import string
79
80
81 def create_batch(papers, batch_size=4):
82     """Create batch of papers for more efficient training"""
83     paper_batches = []
84     for i in range(0, len(papers), batch_size):
85         paper_batches.append(papers[i:i+batch_size])
86     return paper_batches
```

```python
1  import requests
2  import xml.etree.ElementTree as ET
3  from bs4 import BeautifulSoup
4
5
6  def extract_hrefs_from_url_by_title(domains=['cs.AI'], target_title='View HTML'):
7      all_hrefs = []
8      for domain in domains:
9          url = f'https://arxiv.org/list/{domain}/recent?skip=0&show=2000'
10         try:
11             response = requests.get(url)
12             response.raise_for_status()
13             soup = BeautifulSoup(response.text, 'html.parser')
14             for a_tag in soup.find_all('a', title=target_title):
15                 if 'href' in a_tag.attrs:
16                     all_hrefs.append(a_tag['href'])
17         except requests.exceptions.RequestException as e:
18             print(f"Error fetching URL '{url}': {e}")
19     return all_hrefs
20
21
22
23  def fetch_arxiv_ids(domains, max_results=5):
24      """
25      Fetches the ArXiv IDs of papers for the specified domains.
26
27      Args:
28          domains (list): A list of ArXiv subject categories (e.g., ["cs.AI", "physics.hep-th"]).
29          max_results (int): The maximum number of results to fetch per domain (default: 100).
30
31      Returns:
32          list: A list of ArXiv paper IDs.
33      """
34      all_ids = []
35      for domain in domains:
36          url = f"http://export.arxiv.org/api/query?search_query=cat:{domain}&start=0&max_results={max_results}"
37          response = requests.get(url)
38          if response.status_code != 200:
39              print(f"Error fetching data for domain: {domain}")
40              continue
41
42          root = ET.fromstring(response.text)
43          for entry in root.findall("{http://www.w3.org/2005/Atom}entry"):
44              # The ArXiv ID is typically found in the <id> tag.
```

```
45              arxiv_id_full = entry.find("{http://www.w3.org/2005/Atom}id").text
46              # The ID often looks like 'http://arxiv.org/abs/2304.01234v1'.
47              # We want to extract just '2304.01234v1'.
48              arxiv_id = arxiv_id_full.split('/')[-1]
49              all_ids.append(arxiv_id)
50      return all_ids
51
52
53 # Download PDF
54 def download_pdf(pdf_url, save_path="paper.pdf"):
55      response = requests.get(pdf_url)
56      if response.status_code == 200:
57          with open(save_path, "wb") as f:
58              f.write(response.content)
59          return save_path
60      return None
61
62 def extract_main_paper_from_html(html_content):
63      soup = BeautifulSoup(html_content, 'html.parser')
64      main_section = soup.find(id='S1')
65
66      if not main_section:
67          return ""
68
69      # Capture everything starting from the main section
70      main_paper_parts = []
71      current = main_section
72      while current:
73          main_paper_parts.append(str(current))
74          current = current.find_next_sibling()
75
76      return "\n".join(main_paper_parts)
77
78
79 import re # Import regular expressions for cleaning
80
81 def get_body_by_id(html_content, target_id):
82      """
83      Extracts the *entire* inner HTML content of an element with a specific ID.
84      (Kept for reference, but not used for the new requirement)
85
86      Args:
87          html_content (str): The HTML content to parse.
88          target_id (str): The ID of the HTML element whose body content is to be extracted.
89
90      Returns:
91          str: The raw inner HTML content of the element, or None if the ID is not found.
92              Returns an empty string if the element is found but has no content.
```

```
 93
 94        Raises:
 95            TypeError: If html_content is not a string.
 96            TypeError: If target_id is not a string.
 97        """
 98        if not isinstance(html_content, str):
 99            raise TypeError("html_content must be a string.")
100        if not isinstance(target_id, str):
101            raise TypeError("target_id must be a string.")
102
103        soup = BeautifulSoup(html_content, 'html.parser')
104        element = soup.find(id=target_id)  # Find the element by its ID
105
106        if element:
107            return str(element.decode_contents())  # Return the raw inner HTML
108        else:
109            return None  # Return None if the element with the ID is not found
110
111 def extract_paragraph_text(html_content):
112        """
113        Extracts and cleans text content specifically from <p class="ltx_p"> tags within HTML.
114
115        It ignores headings, links, citations, and other non-paragraph elements.
116        It also cleans up citation markers like '[26, 11]' and extra whitespace.
117
118        Args:
119            html_content (str): The HTML content to parse.
120
121        Returns:
122            str: A single string containing the concatenated and cleaned text
123                 from all found <p class="ltx_p"> tags, separated by newlines.
124                 Returns an empty string if no such paragraphs are found.
125
126        Raises:
127            TypeError: If html_content is not a string.
128        """
129        if not isinstance(html_content, str):
130            raise TypeError("html_content must be a string.")
131
132        soup = BeautifulSoup(html_content, 'html.parser')
133        paragraphs = soup.find_all('p', class_='ltx_p') # Find all <p> tags with class 'ltx_p'
134
135        extracted_texts = []
136        for p in paragraphs:
137            # Get text, stripping inner tags like <a>, <cite>, <em>
138            text = p.get_text(separator=' ', strip=True)
139
140            # Use regex to remove citation markers like [26, 11] or [ 23 ]
```

```
141        text = re.sub(r'\[\s*(\d+\s*,?\s*)+\]', '', text)
142
143        # Optional: Clean up potential multiple spaces resulting from tag removal
144        text = re.sub(r'\s+', ' ', text).strip()
145
146        if text: # Add non-empty paragraphs
147            extracted_texts.append(text)
148
149    # Join the texts from all paragraphs with a newline for readability
150    return "\n".join(extracted_texts)
```

```
1 import torch
2
3 def save_checkpoint(model, optimizer, epoch, loss, path="checkpoint.pt"):
4     torch.save({
5         'epoch': epoch,
6         'model_state_dict': model.state_dict(),
7         'optimizer_state_dict': optimizer.state_dict(),
8         'loss': loss
9     }, path)
10    print(f"✅ Checkpoint saved at epoch {epoch} to {path}")
11
12
13 def load_checkpoint(model, optimizer, path="checkpoint.pt"):
14    checkpoint = torch.load(path, map_location=torch.device('cuda' if torch.cuda.is_available() else 'cpu'))
15    model.load_state_dict(checkpoint['model_state_dict'])
16    optimizer.load_state_dict(checkpoint['optimizer_state_dict'])
17    print(f"📦 Loaded checkpoint from epoch {checkpoint['epoch']} with loss {checkpoint['loss']:.4f}")
18    return checkpoint['epoch'], checkpoint['loss']
19
```

```
1 import torch.nn as nn
2 import torch.nn.functional as F
3
4 # First, make sure LuongAttention is defined
5 class LuongAttention(nn.Module):
6     def __init__(self, hidden_dim):
7         super(LuongAttention, self).__init__()
8         self.attn = nn.Linear(hidden_dim, hidden_dim)
9
10    def forward(self, decoder_hidden, encoder_outputs):
11        # decoder_hidden: (batch, hidden)
12        # encoder_outputs: (batch, seq_len, hidden)
13
14        # Transform decoder hidden to match encoder dimension
15        query = self.attn(decoder_hidden).unsqueeze(2)  # (batch, hidden, 1)
16
```

```
17          # Compute scores (dot product)
18          attn_scores = torch.bmm(encoder_outputs, query).squeeze(2)  # (batch, seq_len)
19
20          # Softmax over time dimension
21          attn_weights = F.softmax(attn_scores, dim=1)  # (batch, seq_len)
22
23          # Weighted sum of encoder outputs
24          context = torch.bmm(attn_weights.unsqueeze(1), encoder_outputs)  # (batch, 1, hidden)
25          context = context.squeeze(1)  # (batch, hidden)
26
27          return context, attn_weights
28
29
30
31
32 class Seq2Seq(nn.Module):
33     def __init__(self, encoder, decoder, pad_token_id):
34         super(Seq2Seq, self).__init__()
35         self.encoder = encoder
36         self.decoder = decoder
37         self.pad_token_id = pad_token_id
38
39     def forward(self, src_input_ids, src_attention_mask, tgt_input_ids):
40         encoder_outputs, (hidden, cell) = self.encoder(src_input_ids, src_attention_mask)
41         output = self.decoder(tgt_input_ids, hidden, cell, encoder_outputs)
42         return output
43
44     def generate(self, src_input_ids, src_attention_mask, max_len=100, bos_token_id=None, eos_token_id=None):
45         """Generate sequence for inference"""
46         if bos_token_id is None:
47             bos_token_id = 1  # Default BOS token ID
48         if eos_token_id is None:
49             eos_token_id = 2  # Default EOS token ID
50
51         device = src_input_ids.device
52         batch_size = src_input_ids.size(0)
53
54         # Get encoder outputs
55         encoder_outputs, (hidden, cell) = self.encoder(src_input_ids, src_attention_mask)
56
57         # Initialize decoder input with BOS token
58         decoder_input = torch.tensor([[bos_token_id]] * batch_size, device=device)
59         generated_sequence = [bos_token_id]
60
61         # Generate tokens one by one
62         for _ in range(max_len):
63             # Generate one step
64             next_token_logits, hidden, cell = self.decoder.generate_step(
```

```
65                decoder_input, hidden, cell, encoder_outputs
66            )
67            next_token_id = torch.argmax(next_token_logits, dim=1).item()
68
69            # Stop if EOS token is generated
70            if next_token_id == eos_token_id:
71                generated_sequence.append(next_token_id)
72                break
73
74            generated_sequence.append(next_token_id)
75            decoder_input = torch.tensor([[next_token_id]], device=device)
76
77        return generated_sequence
78
79
80 from torch.nn.utils.rnn import pad_sequence
81
82 import random
83
84 def split_papers(papers, train_ratio=0.8, shuffle=True):
85     """
86     Splits a list of paper URLs into training and testing sets.
87
88     Args:
89         papers (list): List of paper URLs.
90         train_ratio (float): Ratio of training papers.
91         shuffle (bool): Whether to shuffle the papers before splitting.
92
93     Returns:
94         (train_paper, test_paper): Tuple of two lists.
95     """
96     if shuffle:
97         random.shuffle(papers)
98
99     split_index = int(len(papers) * train_ratio)
100     train_paper = papers[:split_index]
101     test_paper = papers[split_index:]
102
103     return train_paper, test_paper
104
```

```
1 !pip install tokenizers
```

```
Requirement already satisfied: tokenizers in /usr/local/lib/python3.11/dist-packages (0.21.1)
Requirement already satisfied: huggingface-hub<1.0,>=0.16.4 in /usr/local/lib/python3.11/dist-packages (from tokenizers) (0.30.1)
Requirement already satisfied: filelock in /usr/local/lib/python3.11/dist-packages (from huggingface-hub<1.0,>=0.16.4->tokenizers) (3.18.0)
Requirement already satisfied: fsspec>=2023.5.0 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub<1.0,>=0.16.4->tokenizers) (2025.3.2)
Requirement already satisfied: packaging>=20.9 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub<1.0,>=0.16.4->tokenizers) (24.2)
```

```
       Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub<1.0,>=0.16.4->tokenizers) (6.0.2)
       Requirement already satisfied: requests in /usr/local/lib/python3.11/dist-packages (from huggingface-hub<1.0,>=0.16.4->tokenizers) (2.32.3)
       Requirement already satisfied: tqdm>=4.42.1 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub<1.0,>=0.16.4->tokenizers) (4.67.1)
       Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub<1.0,>=0.16.4->tokenizers) (4.13.1)
       Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests->huggingface-hub<1.0,>=0.16.4->tokenizers) (3.4.1)
       Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests->huggingface-hub<1.0,>=0.16.4->tokenizers) (3.10)
       Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests->huggingface-hub<1.0,>=0.16.4->tokenizers) (2.3.0)
       Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests->huggingface-hub<1.0,>=0.16.4->tokenizers) (2025.1.31)
```

```python
 1 # papers = extract_hrefs_from_url_by_title()
 2 # papers = papers[1:2]
 3
 4 # for link in papers:
 5 #     save_as = 'paper.html'
 6 #     download_arxiv_html(link, save_as)
 7
 8 #     with open(save_as, 'r', encoding='utf-8') as file:
 9 #         html_data = file.read()
10
11 #     abstract, html_without_abstract = extract_abstract(html_data)
12 #     main_paper = extract_main_paper_from_html(html_without_abstract)
13 #     cleaned_abstract = preprocess_paper_text(abstract)['text']
14 #     cleaned_paper = preprocess_paper_text(main_paper)['text']
15 #     print(cleaned_paper[1000:2000])
```

```python
 1 MODEL_DIR = "Downloads/NLP_Local"
 2 CHECKPOINT_PATH = f"{MODEL_DIR}/checkpoint.pt"
 3 TOKENIZER_PATH = f"{MODEL_DIR}/mytokenizer"
 4 PICKLE_PATH = f"{MODEL_DIR}/model.pkl"
 5
 6 import os
 7 os.makedirs(MODEL_DIR, exist_ok=True)
 8
 9
10 def load_checkpoint(model, optimizer, path="checkpoint.pt"):
11     checkpoint = torch.load(path, map_location=torch.device('cuda' if torch.cuda.is_available() else 'cpu'))
12
13     model.load_state_dict(checkpoint['model_state_dict'])
14     optimizer.load_state_dict(checkpoint['optimizer_state_dict'])
15     epoch = checkpoint['epoch']
16     loss = checkpoint['loss']
17
18     print(f"✅ Loaded checkpoint from epoch {epoch} with loss {loss:.4f}")
19     return epoch, loss
```

```python
 1 import spacy
 2 import torch
```

```python
  3 import torch.nn as nn
  4 import numpy as np
  5
  6 # Load scispaCy model - you'll need to install it first with:
  7 # pip install scispacy
  8 # pip install https://s3-us-west-2.amazonaws.com/ai2-s2-scispacy/releases/v0.5.1/en_core_sci_md-0.5.1.tar.gz
  9 try:
 10     nlp = spacy.load("en_core_sci_md")
 11     nlp.max_length = 2000000
 12 except OSError:
 13     print("Please install the scispaCy model with:")
 14     print("pip install scispacy")
 15     print("pip install https://s3-us-west-2.amazonaws.com/ai2-s2-scispacy/releases/v0.5.1/en_core_sci_md-0.5.1.tar.gz")
 16     raise
 17
 18 from transformers import AutoTokenizer
 19
 20 def tokenize_text(text):
 21     doc = nlp(text)
 22     tokens = [token.text for token in doc if not token.is_stop and token.is_alpha]
 23     return tokens
 24
 25 custom_tokenizer = AutoTokenizer.from_pretrained("allenai/scibert_scivocab_uncased")
 26
 27 custom_tokenizer.bos_token = custom_tokenizer.cls_token  # Use [CLS] as BOS
 28 custom_tokenizer.eos_token = custom_tokenizer.sep_token
 29
 30 class ScispacyEncoder(nn.Module):
 31     def __init__(self, embedding_dim, hidden_dim, num_layers=1, max_sequence_length=50):
 32         super(ScispacyEncoder, self).__init__()
 33         # ScispaCy model has 200-dim embeddings
 34         self.embedding_dim = 200
 35         self.max_sequence_length = max_sequence_length
 36         self.projection = nn.Linear(self.embedding_dim, embedding_dim)
 37         self.lstm = nn.LSTM(embedding_dim, hidden_dim, num_layers, batch_first=True)
 38
 39     def get_embeddings_sequence(self, text):
 40         """Get ScispaCy embeddings for text as a sequence"""
 41         doc = nlp(text)
 42
 43         # Extract important sentences to create a meaningful sequence
 44         # Use basic frequency-based approach to identify key sentences
 45         from collections import Counter
 46
 47         # Count word frequencies (excluding stop words)
 48         word_freq = Counter([token.text.lower() for token in doc
 49                             if not token.is_stop and not token.is_punct
 50                             and token.has_vector])
```

```
51
52             # Score sentences by sum of word frequencies
53             sentences = list(doc.sents)
54             sentence_scores = []
55             for sent in sentences:
56                 score = sum(word_freq[token.text.lower()] for token in sent
57                             if token.has_vector and not token.is_stop and not token.is_punct)
58                 sentence_scores.append((sent, score))
59
60             # Take top sentences up to max_sequence_length
61             top_sentences = sorted(sentence_scores, key=lambda x: x[1], reverse=True)[:self.max_sequence_length]
62             # Re-sort to preserve original order
63             top_sentences = sorted(top_sentences, key=lambda x: sentences.index(x[0]))
64
65             # Get vector for each sentence
66             sequence_vectors = []
67             for sent, _ in top_sentences:
68                 vectors = [token.vector for token in sent if token.has_vector]
69                 if vectors:
70                     mean_vector = np.mean(vectors, axis=0)
71                     # Ensure vector has correct dimensions
72                     if mean_vector.shape[0] != self.embedding_dim:
73                         if mean_vector.shape[0] > self.embedding_dim:
74                             mean_vector = mean_vector[:self.embedding_dim]  # Truncate
75                         else:
76                             # Pad with zeros
77                             padded = np.zeros(self.embedding_dim)
78                             padded[:mean_vector.shape[0]] = mean_vector
79                             mean_vector = padded
80                     sequence_vectors.append(mean_vector)
81                 else:
82                     # Use zeros for sentences with no valid vectors
83                     sequence_vectors.append(np.zeros(self.embedding_dim))
84
85         # Pad or truncate sequence to match max_sequence_length
86         if len(sequence_vectors) > self.max_sequence_length:
87             sequence_vectors = sequence_vectors[:self.max_sequence_length]
88         elif len(sequence_vectors) < self.max_sequence_length:
89             padding_needed = self.max_sequence_length - len(sequence_vectors)
90             for _ in range(padding_needed):
91                 sequence_vectors.append(np.zeros(self.embedding_dim))
92
93         return torch.tensor(np.array(sequence_vectors), dtype=torch.float)
94
95     def forward(self, texts, attention_mask=None):
96         batch_size = len(texts)
97         embedded_sequences = []
98
```

```
 99            # Process each text in the batch
100            for text in texts:
101                # Get sequence of embeddings from ScispaCy
102                seq_embedding = self.get_embeddings_sequence(text)  # [seq_len, embed_dim]
103
104                # Project each vector to desired embedding dimension
105                projected_seq = self.projection(seq_embedding)  # [seq_len, embed_dim]
106                embedded_sequences.append(projected_seq)
107
108            # Stack embeddings
109            embedded = torch.stack(embedded_sequences, dim=0)  # [batch, seq_len, embed_dim]
110
111            # Process through LSTM
112            outputs, (h, c) = self.lstm(embedded)
113
114            return outputs, (h, c)
115
116 class ScispacyDecoder(nn.Module):
117     def __init__(self, vocab_size, embedding_dim, hidden_dim, num_layers=1):
118         super(ScispacyDecoder, self).__init__()
119         self.embedding = nn.Embedding(vocab_size, embedding_dim)
120         self.lstm = nn.LSTM(embedding_dim, hidden_dim, num_layers, batch_first=True)
121         self.attention = LuongAttention(hidden_dim)  # Add attention mechanism
122         # Combine context and hidden for output
123         self.fc_out = nn.Linear(hidden_dim * 2, hidden_dim)
124         self.output_layer = nn.Linear(hidden_dim, vocab_size)
125
126     def forward(self, tgt_input_ids, hidden, cell, encoder_outputs):
127         # [B, L] -> [B, L, D]
128         embedded = self.embedding(tgt_input_ids)
129
130         # Pass through LSTM
131         outputs, (hidden, cell) = self.lstm(embedded, (hidden, cell))
132
133         # Apply attention for each timestep
134         batch_size, seq_len, _ = outputs.size()
135         attention_outputs = []
136
137         for t in range(seq_len):
138             # Get decoder hidden state at this timestep
139             decoder_hidden = outputs[:, t, :]
140
141             # Calculate attention context
142             context, _ = self.attention(decoder_hidden, encoder_outputs)
143
144             # Combine context and hidden state
145             concat_input = torch.cat((decoder_hidden, context), dim=1)
146             output = self.fc_out(concat_input)
```

```
147                    attention_outputs.append(output)
148
149            # Stack attention outputs
150            attention_outputs = torch.stack(attention_outputs, dim=1)
151
152            # Get logits
153            logits = self.output_layer(attention_outputs)
154            return logits
155
156     def generate_step(self, decoder_input, hidden, cell, encoder_outputs):
157            # [1, 1] -> [1, 1, D]
158            embedded = self.embedding(decoder_input)
159
160            # Pass through LSTM for one step
161            outputs, (hidden, cell) = self.lstm(embedded, (hidden, cell))
162
163            # Apply attention
164            decoder_hidden = outputs[:, -1, :]
165            context, _ = self.attention(decoder_hidden, encoder_outputs)
166
167            # Combine context and hidden state
168            concat_input = torch.cat((decoder_hidden, context), dim=1)
169            attention_output = self.fc_out(concat_input)
170
171            # Get logits for the next token
172            logits = self.output_layer(attention_output)
173            return logits, hidden, cell
174
175
176 # Updated encoder-decoder architecture
177 class SciSummarizationModel(nn.Module):
178     def __init__(self, vocab_size, embedding_dim, hidden_dim, num_layers=1):
179            super(SciSummarizationModel, self).__init__()
180            self.encoder = ScispacyEncoder(embedding_dim, hidden_dim, num_layers)
181            # Use the improved decoder with attention
182            self.decoder = ScispacyDecoder(vocab_size, embedding_dim, hidden_dim, num_layers)
183            self.pad_token_id = custom_tokenizer.pad_token_id
184
185     def forward(self, source_texts, tgt_input_ids):
186            encoder_outputs, (hidden, cell) = self.encoder(source_texts)
187            output = self.decoder(tgt_input_ids, hidden, cell, encoder_outputs)
188            return output
189
190     def generate(self, source_text, max_len=100, bos_token_id=None, eos_token_id=None):
191            """Generate sequence for inference"""
192            device = next(self.parameters()).device
193
194            # Use tokenizer's CLS/SEP tokens if BOS/EOS are not available
```

```python
195         if bos_token_id is None:
196             bos_token_id = custom_tokenizer.cls_token_id  # [CLS] token in BERT
197         if eos_token_id is None:
198             eos_token_id = custom_tokenizer.sep_token_id  # [SEP] token in BERT
199
200         # Get encoder outputs
201         encoder_outputs, (hidden, cell) = self.encoder([source_text])
202
203         # Initialize decoder input with BOS token
204         decoder_input = torch.tensor([[bos_token_id]], device=device)
205         generated_sequence = [bos_token_id]
206
207         # Generate tokens one by one
208         for _ in range(max_len):
209             # Use the modified decoder.generate_step method
210             next_token_logits, hidden, cell = self.decoder.generate_step(
211                 decoder_input, hidden, cell, encoder_outputs
212             )
213             next_token_id = torch.argmax(next_token_logits, dim=1).item()
214
215             # Stop if EOS token is generated
216             if next_token_id == eos_token_id:
217                 generated_sequence.append(next_token_id)
218                 break
219
220             generated_sequence.append(next_token_id)
221             decoder_input = torch.tensor([[next_token_id]], device=device)
222
223         return generated_sequence
224
225
226 # Function to preprocess text using ScispaCy
227 import re
228
229 def preprocess_with_scispacy(text):
230     # Remove URLs
231     text = re.sub(r'http\S+|www\.\S+', '', text)
232
233     doc = nlp(text)
234     cleaned_tokens = []
235
236     for token in doc:
237         token_text = token.text
238
239         # Skip stopwords or punctuations
240         if token.is_stop or token.is_punct:
241             continue
242
```

```
243            # Remove wrapped in {}, [], ()
244            if re.match(r'^[\[\(\{].*[\]\)\}]$', token_text):
245                continue
246
247            # Remove tokens with slashes or backslashes
248            if '/' in token_text or '\\' in token_text:
249                continue
250
251            # Remove tokens that contain non-alphanumeric characters
252            if not token_text.isalnum():
253                continue
254
255            # Append cleaned, lemmatized lowercase word
256            cleaned_tokens.append(token.lemma_.lower())
257
258    return " ".join(cleaned_tokens)
259
260
261
262
263
264 # Updated summarize_paper function
265 def summarize_paper_with_scispacy(model, tokenizer, link, max_summary_len=100):
266    device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
267    model.to(device)
268    model.eval()
269
270    # Download and preprocess the paper
271    save_as = 'test_paper.html'
272    download_arxiv_html(link, save_as)
273
274    with open(save_as, 'r', encoding='utf-8') as file:
275        html_data = file.read()
276
277    abstract, html_without_abstract = extract_abstract(html_data)
278    main_paper = extract_paragraph_text(html_without_abstract)
279
280    # Preprocess with ScispaCy
281    cleaned_paper = preprocess_with_scispacy(main_paper)
282    print(f"Clean paper input: {cleaned_paper}")
283    # Generate summary
284    with torch.no_grad():
285        generated_ids = model.generate(cleaned_paper, max_len=max_summary_len)
286
287    # Decode
288    valid_ids = [token_id for token_id in generated_ids if token_id < tokenizer.vocab_size]
289    summary_text = tokenizer.decode(valid_ids, skip_special_tokens=True,clean_up_tokenization_spaces=True)
290
```

```
291     print(" 📄 Original Paper Length:", len(main_paper.split()))
292     print(" 📝 Generated Abstract:", summary_text)
293     print("-" * 60)
294
295     return summary_text
```

```
1 # Model parameters
2 import torch.nn as nn
3 loss_fn = nn.CrossEntropyLoss()
4 vocab_size = len(custom_tokenizer.vocab)
5 embedding_dim = 256  # Increased from 128
6 hidden_dim = 256      # Increased from 128
7 num_epochs = 1
8
9 papers = extract_hrefs_from_url_by_title()
10 papers = papers[1:10]
11 train_paper,test_paper = split_papers(papers)
12 device = torch.device('cpu')
13 # Create model
14
15 special_tokens_dict = {'bos_token': '<s>', 'eos_token': '</s>'}
16 num_added_toks = custom_tokenizer.add_special_tokens(special_tokens_dict)
17
18 # Resize model embeddings
19 sci_model = SciSummarizationModel(vocab_size, embedding_dim, hidden_dim).to(device)
20 optimizer = torch.optim.Adam(sci_model.parameters(), lr=3e-4)
21
22 # Train the model (simplified example)
23 for epoch in range(num_epochs):
24     for link in train_paper:
25         # Get paper data (same as before)
26         save_as = 'paper.html'
27         download_arxiv_html(link, save_as)
28
29         with open(save_as, 'r', encoding='utf-8') as file:
30             html_data = file.read()
31         abstract, html_without_abstract = extract_abstract(html_data)
32         main_paper = extract_paragraph_text(html_without_abstract)
33
34         # Use ScispaCy preprocessing
35         cleaned_paper = preprocess_with_scispacy(main_paper)
36
37         print(f"abstract {abstract}")
38         print("==========================================")
39         # Tokenize abstract for target
40         encoded_abstract = custom_tokenizer(
41             abstract,
42             padding='max_length',
```

```
43              truncation=True,
44              max_length=128,
45              return_tensors='pt',
46              add_special_tokens=True
47          ).to(device)
48
49          # Tokenize cleaned_paper for source input
50          print(f"Inout clean {cleaned_paper}")
51
52          # Decoder input/output setup
53          decoder_input = encoded_abstract.input_ids[:, :-1]  # exclude last token
54          target_labels = encoded_abstract.input_ids[:, 1:]   # exclude first token
55
56          # Forward pass with paper text directly
57          output_logits = sci_model([cleaned_paper], decoder_input)
58
59          # Calculate loss
60          loss = loss_fn(output_logits.view(-1, vocab_size), target_labels.view(-1))
61
62          # Backward pass
63          optimizer.zero_grad()
64          loss.backward()
65          torch.nn.utils.clip_grad_norm_(sci_model.parameters(), 1.0)
66          optimizer.step()
67
68          print(f"Epoch {epoch+1}, Loss: {loss.item():.4f}")
```

abstract We propose the Dual Engines of Thoughts (DEoT), an analytical framework for comprehensive open-ended reasoning. While traditional reasoning frameworks primari
============================================
Inout clean keywords dual engines thoughts analysis framework reasoning framework today interconnected world analyze implication complex event require nuanced grasp in
Epoch 1, Loss: 10.3460
abstract Large Language Models (LLMs) demonstrate impressive capabilities in natural language processing but suffer from inaccuracies and logical inconsistencies known
============================================
Inout clean keyword llm ontology reasoning consistency checking knowledge representation hallucination mitigation hybrid machine learning logical formalism large langu
Epoch 1, Loss: 10.3557
abstract We demonstrate how AI agents can coordinate to deceive oversight systems using automated interpretability of neural networks.
Using sparse autoencoders (SAEs) as our experimental framework, we show that language models (Llama, DeepSeek R1, and Claude 3.7 Sonnet) can generate deceptive explana
Our agents employ steganographic methods to hide information in seemingly innocent explanations, successfully fooling oversight models while achieving explanation qual
We further find that models can scheme to develop deceptive strategies when they believe the detection of harmful features might lead to negative consequences for them
All tested LLM agents were capable of deceiving the overseer while achieving high interpretability scores comparable to those of reference labels.
We conclude by proposing mitigation strategies, emphasizing the critical need for robust understanding and defenses against deception.
============================================
Inout clean sparse autoencoder sae neural network large number neuron use sparsity constraint training call autoencoder approximate identity function ng et 2011 contai
Epoch 1, Loss: 10.3406
abstract Aligning large language models with human preferences is crucial for their safe deployment. While Direct Preference Optimization (DPO) offers an efficient alt
============================================
Inout clean align large language models llm carefully curate human feedback prove critical steer behavior helpful honest harmless response preference optimization meth
Epoch 1, Loss: 10.3329
abstract A popular approach to neurosymbolic AI is to take the output of the last layer of a neural network, e.g. a softmax activation, and pass it through a sparse co
This induces a probability distribution over a set of random variables, which happen to be conditionally independent of each other in many commonly used neurosymbolic

Such conditionally independent random variables have been deemed harmful as their presence has been observed to co-occur with a phenomenon dubbeddeterministic bias, wh
We provide evidence contesting this conclusion and show that the phenomenon ofdeterministic biasis an artifact of improperly applying neurosymbolic AI.
=========================================
Inout clean neurosymbolic nesy ai approach ai seek combine logic neural network integration symbolic method allow inter alia interpretable datum efficient ai system po
Epoch 1, Loss: 10.3469
abstract AlphaZero in 2017 was able to master chess and other games without human knowledge by playing millions of games against itself (self-play), with a computation
=========================================
Inout clean conquer chess holy grail testbe ai development inception supercomputer deep blue ai system beat world champion chess classical time control development har
Epoch 1, Loss: 10.3319
abstract Generative AI is transforming computing education by enabling the automatic generation of personalized content and feedback. We investigate its capabilities i
=========================================
Inout clean generative ai transform learning teaching compute education advanced generative model openai github copilot reshape student teacher experience student mode
Epoch 1, Loss: 10.3236

```
1 # Test your model
2 for link in test_paper[:3]:
3     summary = summarize_paper_with_scispacy(sci_model, custom_tokenizer, link)
```

ean paper input: linecolor gray topline false bottomline false leftline true rightline false backgroundcolor giovanni mauro 1 moruzzi 1 pisa 56124 italy 2 scuola normal
Original Paper Length: 10641
Generated Abstract: pancreatic reaction attractiveness academic commission diesel its its formulations heterologous professorulin transientmentationlmife initiation},1
----------------------------------------------------------
ean paper input: scheduling problem exist dynamic environment unpredictable event unforeseen machine failure arrival urgent job date alteration unexpected weather chang
Original Paper Length: 11077
Generated Abstract: pancreatic reaction attractiveness academic commission diesel its its formulations heterologous professorulin transientmentationlmife initiation},1
----------------------------------------------------------

1 Start coding or generate with AI.