

# 350 NLP Projects

with Code

The Most Powerful NLP-Weapon Arsenal

**Himanshu Ramchandani**  
**M.Tech | Data Science**

## NLP Migrant Workers' Paradise: Almost the most complete Chinese NLP resource library

In the process of getting started and getting familiar with NLP, I used a lot of packages on github, so I sorted it out and shared it here.

Many bags are very interesting and worth collecting, satisfying everyone's collection addiction! If you find it useful, please share and star★,thanks!

Long-term irregular updates, welcome to watch and fork! ❤️❤️❤️



- |                                  |                              |
|----------------------------------|------------------------------|
| * Corpus                         | * Document Processing        |
| * Thesaurus and lexical tools    | * Table Processing           |
| * Pre-trained language model     | * Text Matching              |
| * Extraction                     | * Text Data Enhancement      |
| * Knowledge map                  | * Text Retrieval             |
| * Text generation                | * Reading Comprehension      |
| * Text summarization             | * Sentiment Analysis         |
| * Intelligent question answering | * Common Regular Expressions |
| * Text error correction          | * Speech Processing          |

- \* [Common regular expressions](#)
- \* [Text visualization](#)
- \* [Event extraction](#)
- \* [Text annotation tool](#)
- \* [Machine translation](#)
- \* [Comprehensive tool](#)
- \* [Digital transformation](#)
- \* [Funny and funny tool](#)
- \* [Anaphora resolution](#)
- \* [Course report interview, etc.](#)
- \* [Text clustering](#)
- \* [Competition](#)
- \* [Text classification](#)
- \* [Financial NLP](#)
- \* [Knowledge reasoning](#)
- \* [Medical NLP](#)
- \* [Explainable NLP](#)
- \* [Legal NLP](#)
- \* [Text adversarial attack](#)
- \* [Text generation image](#)
- \* [Others](#)

## corpus

Resource name (Name)	Description	Link
Corpus of names		<a href="#">wainshine/Chinese-Names-Corpus</a>
Chinese-Word-Vectors	Various Chinese word vectors	<a href="#">github repo</a>

Chinese Chat Corpus	The library includes Douban Duolun, PTT gossip corpus, Qingyun corpus, TV drama dialogue corpus, Tieba forum reply corpus, Weibo corpus, little yellow chicken corpus	<a href="#">link</a>
Chinese rumor data	In this data file, each line contains a rumor data in json format	<a href="#">github</a>
Chinese Question Answering Dataset		<a href="#">link</a> extract code 2dva
WeChat official account corpus	The 3G corpus, which includes some articles from WeChat official accounts captured from the web, has removed HTML and only contains plain text. One article per line, in JSON format, name is the name of the WeChat official account, account is the ID of the WeChat official account, title is the title, and content is the text	<a href="#">github</a>
Chinese natural language processing corpus, data set		<a href="#">github</a>

Task-based dialogue English dataset	[The most complete task-based dialogue data set] mainly introduces a complete task-based dialogue data set, which covers the main information of all commonly used data sets in the field of task-based dialogue. In addition, in order to help researchers better grasp the context of field progress, we present the State-of-the-art experimental results on several datasets in the form of Leaderboard.	<a href="#">github</a>
Speech Recognition Corpus Generation Tool	Create an Automatic Speech Recognition (ASR) corpus from online videos with audio/subtitles	<a href="#">github</a>
LitBank NLP dataset	A corpus of 100 labeled English novels supporting natural language processing and computational humanities tasks	<a href="#">github</a>
ChineseULMFiT	Sentiment Analysis Text Classification Corpus and Model	<a href="#">github</a>
The administrative division data of provinces, municipalities and towns are marked with pinyin		<a href="#">github</a>

Automated Summarization Corpus of Education Industry News		<a href="#">github</a>
Chinese Natural Language Processing Dataset		<a href="#">github</a>
Baidu Zhizhi Q&A Corpus	More than 5.8 million questions, 9.38 million answers, 5800 classification labels. Based on the question and answer corpus, it can support a variety of applications, such as chat question and answer, logic mining	<a href="#">github</a>
Wikipedia Massively Parallel Text Corpus	85 languages, 1620 language pairs, 135M contrasting sentences	<a href="#">github</a>
Ancient Poetry Thesaurus		<a href="#">github repo</a>  <a href="#">more complete ancient poetry lexicon</a>
Low memory loading Wikipedia data	Use the new version of nlp library to load 17GB+ English Wiki corpus and only occupy 9MB of memory Traversal speed 2-3 Gbit/s	<a href="#">github</a>
couplet data	700,000 couplets, more than 700,000 couplets	<a href="#">github</a>

"Color Dictionary" dataset		<a href="#">github</a>
42GB of JD Customer Service Dialogue Data (CSDD)		<a href="#">github</a>
700,000 couplet data		<a href="#">link</a>
Username Blacklist List		<a href="#">github</a>
Dependency parsing corpus	40,000 high-quality labeled data	<a href="#">Homepage</a>
People's Daily Corpus Processing Toolset		<a href="#">github</a>
False news dataset fake news corpus		<a href="#">github</a>
Poetry Quality Evaluation / Fine-grained Emotional Poetry Corpus		<a href="#">github</a>
Open tasks related to Chinese natural language processing	Dataset and current best results	<a href="#">github</a>
Chinese abbreviation dataset		<a href="#">github</a>

Chinese task benchmarking	Representative dataset - benchmark (pretrained) model - corpus - baseline - toolkit - leaderboard	<a href="#">github</a>
Chinese Rumor Database		<a href="#">github</a>
CLUEDatasetSearch	Chinese and English NLP datasets Search all Chinese NLP datasets, with commonly used English NLP datasets attached	<a href="#">github</a>
Multi-Document Summarization Dataset		<a href="#">github</a>
Make Everyone "Courteous" Polite Migration Quest	Transform impolite sentences into polite ones while preserving meaning, providing a dataset with 139M+ instances	<a href="#">paper and code</a>
Cantonese/English Conversational Bilingual Corpus		<a href="#">github</a>
List of Chinese NLP datasets		<a href="#">github</a>
Nomenclature recognition data set of person-like names/place		<a href="#">github</a>



names/organization  
names

Chinese Language  
Comprehension  
Benchmark

Includes representative  
datasets & benchmark models  
& corpora & leaderboards

[github](#)

OpenCLaP  
multi-domain open  
source Chinese  
pre-trained language  
model warehouse

Civil documents, criminal  
documents, Baidu  
Encyclopedia

[github](#)

Chinese full word  
coverage BERT and  
two reading  
comprehension data

DRCD dataset: Released by  
Delta Research Institute of  
Taiwan, China, it has the same  
form as SQuAD, and is an  
extractive reading  
comprehension dataset based  
on traditional Chinese.

[github](#)

CMRC 2018 dataset: Chinese  
machine reading  
comprehension data released  
by the Xunfei Joint Laboratory  
of Harbin Institute of  
Technology. According to a  
given question, the system  
needs to extract fragments  
from the text as answers, in  
the same form as SQuAD.

Dakshina dataset

Latin/native script parallel  
dataset for twelve South Asian  
languages

[github](#)

OPUS-100	Multilingual (100 kinds) parallel corpus centered on English	<a href="#">github</a>
Chinese Reading Comprehension Dataset		<a href="#">github</a>
Chinese natural language processing vector collection		<a href="#">github</a>
Chinese Language Comprehension Benchmark	Includes representative datasets, benchmark (pretrained) models, corpora, leaderboards	<a href="#">github</a>
Large list of NLP datasets/benchmark tasks		<a href="#">github</a>
LitBank NLP dataset	A corpus of 100 labeled English novels supporting natural language processing and computational humanities tasks	<a href="#">github</a>
700,000 couplet data		<a href="#">github</a>
Parallel Corpus of Classical Chinese (Ancient Chinese)-Modern Chinese	The short chapters include "The Analects of Confucius", "Mencius", "Zuo Zhuan" and other short ancient books, which have been merged with "Zi Zhi Tong Jian"	<a href="#">github</a>

COLDDateset,  
Chinese Offensive  
Language Detection  
Dataset

Covers topics such as race,  
gender, and region, and the  
data will be released after the  
paper is published

[paper](#)

## Thesaurus and Lexical Tools

Resource name (Name)	Description	Link
textfilter	Sensitive word filtering in Chinese and English	<a href="#">observerss/textfilter</a>
Name extraction function	Chinese (modern, ancient) names, Japanese names, Chinese surnames and first names, titles (big aunt, little aunt, etc.), English -> Chinese name (John Lee), idiom dictionary	<a href="#">cocoNLP</a>
Chinese Abbreviation Library	National People's Congress: National People's Congress; China: People's Republic of China; Women's Tennis: Women/n Tennis/n Game/vn	<a href="#">github</a>

Chinese Dictionaries	How to dismantle Chinese characters (1) How to dismantle (2) How to dismantle (3)	<a href="#">kfcd/chaizi</a>
Lexical Sentiment Value	Mountain spring water: 0.400704566541 Sufficient  : 0.37006739587	<a href="#">rainarch/SentiBridge</a>
Chinese thesaurus, stop words, sensitive words		<a href="#">dongxiexidian/Chinese</a>
python-pinyin	Convert Chinese characters to Pinyin	<a href="#">mozillazg/python-pinyin</a>
zhtools	Conversion between Traditional and Simplified Chinese	<a href="#">skydark/nstools</a>
English simulation Chinese pronunciation engine	say wo i ni #say: I love you	<a href="#">tinyfool/ChineseWithEnglish</a>
chinese_dictionary	Thesaurus, antonym, negative thesaurus	<a href="#">guotong1988/chinese_dictionary</a>

wordninja

English string  
segmentation and  
word extraction  
without spaces

[wordninja](#)

Vocabulary related  
to automobile brand  
and automobile  
parts

[data](#)

Thesaurus  
organized by THU

IT thesaurus,  
financial  
thesaurus, idiom  
thesaurus, place  
names, historical  
celebrity  
thesaurus, poetry  
thesaurus,  
medical  
thesaurus, diet  
thesaurus, legal  
thesaurus,  
automobile  
thesaurus, animal  
thesaurus

[link](#)

Crime Legal Terms  
and Classification  
Model

Contains 856  
crime knowledge  
graphs, crime  
prediction based  
on 2.8 million  
crime training  
database, 13  
types of question  
classification and  
legal information  
question and  
answer function  
based on 20W  
legal question and  
answer pairs

[github](#)

Word segmentation  
corpus + code

[Baidu network disk link](#) -  
extraction code pea6

Chinese word  
segmentation +  
part-of-speech  
tagging based on  
Bi-LSTM + CRF

keras  
implementation

[link](#)

Chinese word  
segmentation and  
part-of-speech  
tagging based on  
Universal  
Transformer + CRF

[link](#)

Fast Neural Network  
Word Segmentation  
Package

java version

chinese-xinhua	Zhonghua Xinhua dictionary database and api, including commonly used Xiehouyu, idioms, words and Chinese characters	<a href="#">github</a>
SpaCy Chinese model	Contains Parser, NER, syntax tree and other functions. Some English packages use spacy's English model. If you want to adapt to Chinese, you may need to use spacy's Chinese model.	<a href="#">github</a>
Chinese character data		<a href="#">github</a>
Synonyms Chinese Synonym Toolkit		<a href="#">github</a>
Harvest Text	Domain adaptive text mining tools (new word discovery-sentiment analysis-entity linking, etc.)	<a href="#">github</a>

word2word	Easy-to-use multilingual word-word pair set 62 languages/3,564 multilingual pairs	<a href="#">github</a>
Polyphone dictionary data and codes		<a href="#">github</a>
Chinese characters, words, idioms query interface		<a href="#">github</a>
103976 English vocabulary packs	(sql version, csv version, Excel version)	<a href="#">github</a>
Big list of swear words in English		<a href="#">github</a>
word pinyin data		<a href="#">github</a>
Number calling library in 186 languages		<a href="#">github</a>
Large-scale name database of countries around the world		<a href="#">github</a>



Chinese character  
feature extractor  
(featurizer)

Extract the  
features of  
Chinese  
characters  
(pronunciation  
features, font  
features) for deep  
learning features

[github](#)

char\_featurizer -  
Chinese character  
feature extraction  
tool

[github](#)

Python interface  
library of mecab, the  
CJK word  
segmentation library

[github](#)

g2pC context-based  
Chinese  
pronunciation  
automatic marking  
module

[github](#)

ssc, Sound Shape  
Code

Phonetic code -  
Chinese character  
string similarity  
calculation method  
based on  
"phonetic code"

[version 1](#)

[version 2](#)

[blog/introduction](#)

Acquisition of multiple meanings/sense items of Chinese words and semantic disambiguation of specific sentences based on the encyclopedia knowledge base

[github](#)

Tokenizer is a fast and customizable text tokenization library

[github](#)

Tokenizers

State-of-the-art tokenizer with a focus on performance and versatility

[github](#)

Realize text "face changing" through synonym replacement

[github](#)

token2index is a powerful lightweight term index library compatible with PyTorch/Tensorflow

[github](#)

Traditional and Simplified Conversion

[github](#)

Cantonese NLP Tools		<a href="#">github</a>
domain dictionary	Professional dictionary knowledge base covering 68 fields with a total of 9.16 million words	<a href="#">github</a>

## Pre-trained language model & large model

Resource name (Name)	Description	Link
BMList	Big Model Big List	<a href="#">github</a>
Chinese translation of bert papers		<a href="#">link</a>
The slides of the original author of bert		<a href="#">link</a>
Text Classification Practice		<a href="#">github</a>
bert tutorial text classification tutorial		<a href="#">github</a>
Bert pytorch implementation		<a href="#">github</a>
Bert pytorch implementation		<a href="#">github</a>

BERT generates sentence vectors, BERT does text classification and text similarity calculation

[github](#)

Diagram of bert and ELMO

[github](#)

BERT Pre-trained models and downstream applications

[github](#)

Language/Knowledge Representation Tool BERT & ERNIE

[github](#)

Using the gpt-2 language model in Kashgari

[github](#)

Facebook LAMA

Probes for analyzing factual and commonsense knowledge contained in pretrained language models. Language model analysis, providing a unified access interface for Transformer-XL/BERT/ELMo/GPT pre-trained language models

[github](#)

Chinese GPT2 training code

[github](#)

XLMTFacebook's cross-language pre-trained language model

[github](#)

Massive Chinese pre-trained ALBERT model

[github](#)

Transformers 20	Supports TensorFlow 20 and PyTorch's natural language processing pre-trained language models (BERT, GPT-2, RoBERTa, XLM, DistilBert, XLNet...) 8 architectures/33 pre-trained models/102 languages	<a href="#">github</a>
8 papers sort out the progress and reflection of BERT related models		<a href="#">github</a>
French RoBERTa pre-trained language model	French RoBERTa pre-trained language model trained with 138GB corpus	<a href="#">link</a>
Chinese pre-trained ELECTREA model	Pretrain Chinese Model based on confrontational learning	<a href="#">github</a>
albert-chinese-ner	Use the pre-trained language model ALBERT to do Chinese NER	<a href="#">github</a>
Open source pre-trained language model collection		<a href="#">github</a>
Chinese ELECTRA pre-training model		<a href="#">github</a>
Predicting Next Word with Transformers (BERT, XLNet, Bart, Electra, Roberta, XLM-Roberta) (Model Comparison)		<a href="#">github</a>
TensorFlow Hub	New language models for 40+ languages (including Chinese)	<a href="#">link</a>

UER	Chinese pre-trained model warehouses based on different corpora, encoders, and target tasks (including BERT, GPT, ELMO, etc.)	<a href="#">github</a>
Open source pre-trained language model collection		<a href="#">github</a>
Multilingual sentence vector package		<a href="#">github</a>
Language Model as a Service (LMaaS)	Language Model as a Service	<a href="#">github</a>
Open source language model GPT-NeoX-20B	20 billion parameters, currently the largest publicly accessible pre-trained general autoregressive language model	<a href="#">github</a>
Chinese Science Literature Dataset (CSL)	Contains 396,209 meta-information (titles, abstracts, keywords, disciplines, categories) of papers in Chinese core journals. The CSL dataset can be used as a pre-training corpus, and can also be used to construct many NLP tasks, such as text summarization (title prediction), keyword generation, and text classification.	<a href="#">github</a>
Large model development artifact		<a href="#">github</a>

**extract**

Resource name (Name)	Description	Link
time extraction	It has been integrated into the python package <a href="#">cocoNLP</a> , welcome to try	<a href="#">java version</a> <a href="#">python version</a>
Neural network relationship extraction pytorch	Chinese is not supported yet	<a href="#">github</a>
Bert-based named entity recognition pytorch	Chinese is not supported yet	<a href="#">github</a>
Keyword (Keyphrase) extraction package pke		<a href="#">github</a>
BLINK's most advanced entity link library		<a href="#">github</a>
Named entity recognition implemented by BERT/CRF		<a href="#">github</a>
Support batch parallel LatticeLSTM Chinese named entity recognition		<a href="#">github</a>
Building a Model for Medical Entity Recognition	Contains dictionaries and corpus annotations, based on python	<a href="#">github</a>

Pipeline Entity and Relationship Extraction Based on TensorFlow and BERT	- Entity and Relation Extraction Based on TensorFlow and BERT Pipeline entity and relationship extraction based on TensorFlow and BERT, the solution to the information extraction task of the 2019 Language and Intelligence Technology Competition. Schema based Knowledge Extraction, SKE 2019	<a href="#">github</a>
Chinese named entity recognition NeuroNER vs BertNER		<a href="#">github</a>
Chinese Named Entity Recognition Based on BERT		<a href="#">github</a>
Chinese key phrase extraction tool		<a href="#">github</a>
bert	tensorflow version for Chinese named entity recognition	<a href="#">github</a>
bert-Kashgari	Kashgari, a keras-based encapsulation classification and labeling framework, can build a classification or sequence labeling model in a few minutes	<a href="#">github</a>
cocoNLP	Extraction of information such as name, address, email address, mobile phone number, mobile phone attribution, etc., rake phrase extraction algorithm.	<a href="#">github</a>



Microsoft Multilingual  
Number/Unit/Eg Date Time  
Recognition Package

[github](#)

Baidu open source benchmark  
information extraction system

[github](#)

Chinese address word  
segmentation (identification and  
extraction of address elements),  
NER through sequence  
annotation

[github](#)

Open Domain Text Knowledge  
Triple Extraction and  
Knowledge Base Construction  
Based on Dependency Syntax

[github](#)

Chinese keyword extraction  
method based on pre-training  
model

[github](#)

chinese\_keyphrase\_extractor  
(CKPE)

A tool for chinese keyphrase  
extraction A tool for quickly  
extracting and identifying  
keyphrases from natural  
language text

[github](#)

Simple resume parser to extract  
key information from resumes

[github](#)

BERT-NER-Pytorch three  
different modes of BERT  
Chinese NER experiments

[github](#)

# knowledge map

Resource name (Name)	Description	Link
Tsinghua University XLORE Chinese-English cross-language encyclopedia knowledge map	Baidu, Chinese Wiki, English Wiki	<a href="#">link</a>
Automatic generation of document maps		<a href="#">github</a>
Question answering system based on knowledge graph in medical field		<a href="#">github</a>  This repo refers to <a href="#">github</a>
Chinese character relationship knowledge map project		<a href="#">github</a>
AmpliGraph Knowledge Graph Representation Learning (Python) Library Knowledge Graph Concept Link Prediction		<a href="#">github</a>
Chinese knowledge map materials, data and tools		<a href="#">github</a>
Chinese Knowledge Graph Based on Baidu Encyclopedia	Extract triplet information and build a Chinese knowledge map	<a href="#">github</a>

Zincbase Knowledge  
Graph Construction Toolkit

[github](#)

Question answering  
system based on  
knowledge graph

[github](#)

Collation of knowledge  
map deep learning related  
materials

[github](#)

Southeast University  
"Knowledge Graph"  
graduate course (data)

[github](#)

Knowledge map car audio  
work project

[github](#)

"One Piece" Knowledge  
Graph

[github](#)

A dataset of 132  
knowledge graphs

Covers common sense, city,  
finance, agriculture, geography,  
weather, social networking,  
Internet of Things, medical care,  
entertainment, life, business,  
travel, science and education

[link](#)

Large-scale, structured,  
Chinese-English bilingual  
COVID-19 Knowledge  
Graph (COKG-19)

[link](#)

Event Triple Extraction  
Based on Dependency  
Syntax and Semantic Role  
Labeling

[github](#)

Abstract Knowledge  
Graph

The current scale is 500,000,  
supporting the abstraction of  
nominal entities, state  
descriptions, and event actions

[github](#)

Large-scale Chinese  
knowledge map data 1.4  
billion entities

[github](#)

Jiagu natural language  
processing tool

Based on models such as  
BiLSTM, it provides functions  
such as knowledge graph  
relationship extraction, Chinese  
word segmentation,  
part-of-speech tagging, named  
entity recognition, sentiment  
analysis, new word discovery,  
keyword text summarization,  
text clustering, etc.

[github](#)

medical\_NER - Chinese  
Medical Knowledge Graph  
Named Entity Recognition

[github](#)

A large list of learning  
materials/datasets/tool  
resources related to  
knowledge graphs

[github](#)

LibKGE is a knowledge graph embedding library for reproducible research

[github](#)

Military field knowledge map question answering project based on mongodb storage

Including aircraft, space equipment, etc. 8 categories, more than 100 subcategories, a total of 5,800 items of military weapons knowledge base, the project does not use a graph database for storage, through jieba to analyze questions, identify entity items in questions, and complete based on query templates The query of multiple types of questions is mainly to provide a demo of the question-and-answer thinking in the industry.

[github](#)

Jingdong Commodity Knowledge Graph

[github](#)

Chinese Relation Extraction Based on Distant Supervision

[github](#)

Intelligent Question Answering System Based on Medical Knowledge Graph

[github](#)

BLINK's most advanced entity link library

[github](#)

A small securities knowledge graph/knowledge base		<a href="#">github</a>
dstlr unstructured text scalable knowledge map construction platform		<a href="#">github</a>
Baidu Encyclopedia character entry attribute extraction	Using BERT-based fine-tuning and feature extraction methods for knowledge graphs	<a href="#">github</a>
Data related to COVID-19	New crown and other types of pneumonia Chinese medical dialogue dataset; open data sources of institutions such as Tsinghua University (COVID-19)	<a href="#">github</a> <a href="#">github</a>
DGL-KE Graph Embedding Representation Learning Algorithm		<a href="#">github</a>
causality map		<a href="#">method data</a>
Causal Event Pairs Based on Multi-Domain Text Datasets		<a href="#">link</a>

## text generation

Resource name (Name)	Description	Link
----------------------	-------------	------

Texar	Toolkit for Text Generation and Beyond	<a href="#">github</a>
Prof. Ehud Reiter's Blog		<a href="#">link</a> Professor Wan Xiaojun of Peking University strongly recommends this blog, which conducts in-depth discussions and reflections on NLG technology, evaluation and application.
Large list of resources related to text generation		<a href="#">github</a>
Open Domain Dialogue Generation and Its Practice in Microsoft Xiaoice	Natural language generation allows machines to master the ability of automatic creation	<a href="#">link</a>
Text Generation Control		<a href="#">github</a>
A large list of natural language generation related resources		<a href="#">github</a>
Evaluating Natural Language Generation with BLEURT		<a href="#">link</a>
Automatic couplet data and robots		<a href="#">Code link</a>

		700,000 couplet data
Automatically generate comments	Generating comments based on Hacker News article titles using Transformer codec model	<a href="#">github</a>
Natural language generation SQL statement (English)		<a href="#">github</a>
Natural Language Generation Resource Collection		<a href="#">github</a>
Benchmarking Chinese Generation Tasks		<a href="#">github</a>
Topic-specific text generation/text augmentation based on GPT2		<a href="#">github</a>
Encoding, Tokenization, and Implementation of a Controlled and Efficient Text Generation Methodology		<a href="#">github</a>
TextFooler's adversarial text generation module for text classification/inference		<a href="#">github</a>



SimBERT	BERT model based on UniLM idea, integrating retrieval and generation	<a href="#">github</a>
New word generation and sentence making	Non-existing words generate new words from scratch with GPT-2 variants, their definitions, and example sentences	<a href="#">github</a>
Automatically generate multiple choice questions from text		<a href="#">github</a>
Synthetic Data Generation Benchmark		<a href="#">github</a>

## text summary

Resource name (Name)	Description	Link
Chinese text summarization/keyword extraction		<a href="#">github</a>
Automatic Summarization of Resume Based on Named Entity Recognition		<a href="#">github</a>

Automatic text summarization library TextTeaser	English only	<a href="#">github</a>
Extractive summary extraction based on the latest language models such as BERT		<a href="#">github</a>
A Comprehensive Guide to Text Summarization with Deep Learning in Python		<a href="#">link</a>
(Colab) Abstract Text Summary Implementation Highlights (Tutorial		<a href="#">github</a>

## Smart Q&A

Resource name (Name)	Description	Link
Chinese chatbot	Train the chatbot you want according to your own corpus, which can be used in scenarios such as intelligent customer service, online question and answer, intelligent chat, etc.	<a href="#">github</a>
Interesting robot qingyun	Chinese chatbot trained by qingyun	<a href="#">github</a>
Open dialogue robots, knowledge graphs, semantic understanding, natural language processing tools and data		<a href="#">github</a>

qa right robot	A model-for-Retrieval chatbot - customer service robot, Chinese Retrieval chatbot (Chinese retrieval robot)	<a href="#">git</a>
ConvLab open source multi-domain end-to-end dialogue system platform		<a href="#">github</a>
A dialog system based on the latest version of rasa		<a href="#">github</a>
Chatbots based on the financial-judicial domain (with the nature of small talk)		<a href="#">github</a>
End-to-end closed-domain dialogue system		<a href="#">github</a>
MiningZhiDaoQACorpus	5.8 million Baidu Zhizhi Q&A data mining project, Baidu Zhizhi Q&A corpus, including more than 5.8 million questions, each with a question label. Based on this question and answer corpus, it can support a variety of applications, such as logic mining	<a href="#">github</a>
GPT2 model GPT2-chitchat for Chinese chatting		<a href="#">github</a>
Selection of relevant resource lists (Leaderboards, Datasets, Papers) based on multiple		<a href="#">github</a>

rounds of responses from  
retrieval chatbots

Microsoft Conversational Bot  
Framework

[github](#)

chatbot-list

Application and architecture of  
intelligent customer service and  
chatbots, algorithm sharing and  
introduction in the industry

[github](#)

Chinese medical dialogue data  
Chinese medical dialogue data  
set

[github](#)

A Large-Scale Medical Dialogue  
Dataset

Contains 1.1 million medical  
consultations and 4 million  
doctor-patient dialogues

[github](#)

Large-scale cross-domain  
Chinese task-oriented  
multi-round dialogue dataset  
and model CrossWOZ

[paper  
& data](#)

Open source conversational  
information search platform

[github](#)

Contextual Interaction  
Multimodal Dialogue Challenge  
2020 (DSTC9 2020)

[github](#)

Use Quora questions to  
paraphrase the trained T5  
questions (Paraphrase)

[github](#)

Google releases Taskmaster-2 natural language task dialogue dataset

[github](#)

Haystack's flexible, powerful, and extensible Question Answering (QA) framework

[github](#)

End-to-end closed-domain dialogue system

[github](#)

Amazon releases knowledge-based human-human open domain dialogue dataset

[github](#)

Albert Large QA model trained based on Baidu webqa and dureader dataset

[github](#)

CommonsenseQA Commonsense-Oriented English QA Challenge

[link](#)

MedQuAD (English) Medical Question Answering Dataset

[github](#)

A Q&A engine using Wikipedia text as context, based on Albert and Electra

[github](#)

A question answering attempt based on the 14W song knowledge base

Functions include Lyrics Solitaire, Finding Songs with Known Lyrics, and Questions and Answers about the

[github](#)

## text error correction

Resource name (Name)	Description	Link
Chinese text error correction module code		<a href="#">github</a>
English spell checking library		<a href="#">github</a>
Python spell checking library		<a href="#">github</a>
GitHub Typo Corpus Large-Scale GitHub Multilingual Spelling/Grammar Error Dataset		<a href="#">github</a>
BertPunc BERT-based state-of-the-art punctuation repair model		<a href="#">github</a>
Chinese writing proofreading tool		<a href="#">github</a>
Text Error Correction Literature List	Chinese Spell Checking (CSC) and Grammatical Error Correction (GEC)	<a href="#">github</a>
Winner of Text Smart Proofreading Contest	It has been applied, from the team of Soochow University and Dharma Academy	<a href="#">link</a>

## multimodal

Resource name (Name)	Description	Link
Chinese Multimodal Dataset "Wukong"	Huawei's Noah's Ark Laboratory open source large-scale, including 100 million text pairs	<a href="#">github</a>
Chinese graphic representation pre-training model Chinese-CLIP	The Chinese version of the CLIP pre-training model, open source multiple model scales, and a few lines of code can handle Chinese image-text representation extraction & image-text retrieval	<a href="#">github</a>

## speech processing

Resource name (Name)	Description	Link
ASR Speech Dataset + Chinese Speech Recognition System Based on Deep Learning		<a href="#">github</a>

Tsinghua University  
THCHS30 Chinese Speech  
Dataset

[data\\_thchs30tgz-OpenSLR domestic image](#)

[data\\_thchs30tgz](#)

[test-noisetgz-OpenSLR domestic image test-noisetgz](#)

[resourcetgz-OpenSLR domestic image](#)

[resourcetgz](#)

[Free ST Chinese Mandarin Corpus](#)

[Free ST Chinese Mandarin Corpus](#)

[AIShell-1 open source version dataset-OpenSLR domestic image](#)

[AIShell-1 open source version dataset](#)

[Primewords Chinese Corpus Set 1-OpenSLR Domestic Mirror](#)

[Primewords Chinese Corpus Set 1](#)

[laughter detector](#)

[github](#)



Common Voice Speech Recognition Dataset New Version	Includes over 1,400 hours of speech samples from 42,000 contributors, covering github	<a href="#">link</a>
speech-aligner	A tool for generating phoneme-level time-aligned annotations from "human voice speech" and its "language text"	<a href="#">github</a>
ASR Speech Dictionary/Dictionary		<a href="#">github</a>
Speech Sentiment Analysis		<a href="#">github</a>
masr	Chinese speech recognition, providing pre-training model, high recognition rate	<a href="#">github</a>
Chinese Text Normalization for Speech Recognition		<a href="#">github</a>
Voice quality evaluation indicators (MOSNet, BSSEval, STOI, PESQ, SRMR)		<a href="#">github</a>
Chinese/English Pronunciation Dictionary for Speech Recognition		<a href="#">github</a>

Multilingual speech-text translation corpus released by CoVoSTFacebook

Includes audio, text transcription and English translation in 11 languages (French, German, Dutch, Russian, Spanish, Italian, Turkish, Persian, Swedish, Mongolian and Chinese)

[github](#)

Parakeet text-to-speech synthesis based on PaddlePaddle

[github](#)

(Java) Accurate Speech Natural Language Detection Library

[github](#)

Multilingual speech-text translation corpus released by CoVoSTFacebook

[github](#)

Text-to-Speech Synthesis Implemented in TensorFlow 2

[github](#)

Python audio feature extraction package

[github](#)

ViSQOL audio quality perception is objective and complete reference index, divided into two modes: audio and voice

[github](#)

zhrtvc	Easy-to-use Chinese voice clone and Chinese speech synthesis system	<a href="#">github</a>
aukit	An easy-to-use speech processing toolbox, including speech noise reduction, audio format conversion, feature spectrum generation and other modules	<a href="#">github</a>
phkit	An easy-to-use phoneme processing toolbox, including Chinese phonemes, English phonemes, text-to-pinyin, text regularization and other modules	<a href="#">github</a>
zhvoice	Chinese speech corpus, the speech is clearer and more natural, including 8 open source data sets, 3200 speakers, 900 hours of speech, 13 million words	<a href="#">github</a>
audio for speech behavior detection	, binarization, speaker recognition, automatic speech recognition, emotion recognition and other audio annotation tools	<a href="#">github</a>

Deep Learning Emotional Text-to-Speech Synthesis	<a href="#">github</a>
Python audio data augmentation library	<a href="#">github</a>
Audio Enhancement Based on Large-Scale Audio Dataset Audioset	<a href="#">github</a>
voice transfer	<a href="#">github</a>

# document processing

Resource name (Name)	Description	Link
LayoutLM-v3 Document Understanding Model		<a href="#">github</a>
PyLaia Deep Learning Toolkit for Handwritten Document Analysis		<a href="#">github</a>
Single-document unsupervised keyword extraction		<a href="#">github</a>

DocSearch Free Documentation Search Engine		<a href="#">github</a>
fdfgen	Ability to automatically create pdf documents and fill in information	<a href="#">link</a>
pdfx	Automatically extract cited references and download the corresponding pdf file	<a href="#">link</a>
invoice2data	Invoice pdf information extraction	<a href="#">invoice2data</a>
PDF document information extraction		<a href="#">github</a>
PDFMiner	PDFMiner can get the exact position of the text in the page, as well as other information such as font or line. It also has a PDF converter that can convert PDF files to other text formats such as HTML. There is also an extensible parser PDF that can be used for other purposes than text analysis.	<a href="#">link</a>
PyPDF2	PyPDF 2 is a python PDF library capable of splitting, merging, cropping and converting pages of PDF files. It can also add custom data, viewing options and passwords to PDF files. It can retrieve text and metadata from PDFs, and can also merge entire files together.	<a href="#">link</a>

PyPDF2	PyPDF 2 is a python PDF library capable of splitting, merging, cropping and converting pages of PDF files. It can also add custom data, viewing options and passwords to PDF files. It can retrieve text and metadata from PDFs, and can also merge entire files together.	<a href="#">link</a>
ReportLab	ReportLab can quickly create PDF documents. A time-proven, super-easy-to-use open source project for creating complex, data-driven PDF documents and custom vector graphics. It's free, open source, and written in Python. With more than 50,000 downloads per month, the package is part of standard Linux distributions, embedded in many products, and was chosen to power Wikipedia's print/export functionality.	<a href="#">link</a>
Simple PDF file text editor written by SIMPdfPython		<a href="#">github</a>
pdf-diff	PDF file diff tool can display the difference between two pdf documents	<a href="#">github</a>

## form processing

Resource name (Name)	Description	Link
Use unet to realize automatic detection of		<a href="#">github</a>

## document tables and table reconstruction

pdftabextract	Used for form information analysis after OCR recognition, very powerful	<a href="#">link</a>
tabula-py	Directly convert the table information in pdf to pandas dataframe, there are two versions of codes in java and python	
camelot	PDF form parsing	<a href="#">link</a>
pdfplumber	PDF form parsing	
PubLayNet	Able to divide paragraphs, identify tables, pictures	<a href="#">link</a>
Extract tabular data from papers		<a href="#">github</a>
Finding answers in tables with BERT		<a href="#">github</a>
Series of articles on table questions and answers		<a href="#">Introduction to the end of the model</a>
Generate tabular data using GAN (English only)		<a href="#">github</a>

carefree-learn (PyTorch)	Automated Machine Learning (AutoML) Package for Tabular Datasets	<a href="#">github</a>
Closed domain fine-tuning table detection		<a href="#">github</a>
PDF form data extraction tool		<a href="#">github</a>
TaBERT A New Model for Understanding Tabular Data Queries		<a href="#">paper</a>
form processing	Awesome-Table-Recognition	<a href="#">github</a>

## text match

Resource name (Name)	Description	Link
Sentence, QA similarity matching MatchZoo	A collection of text similarity matching algorithms, including multiple deep learning methods, worth trying.	<a href="#">github</a>
Chinese Question Sentence Similarity Calculation Competition and Scheme Summary		<a href="#">github</a>
similarity similarity calculation toolkit	Written in java, it is used for similarity calculations related to words, phrases, sentences, lexical analysis,	<a href="#">github</a>



	sentiment analysis, semantic analysis, etc.	
Chinese word similarity calculation method	Combined with the word similarity calculation method of Synonyms Cilin Extended Edition and Hownet, the vocabulary coverage is more and the results are more accurate.	<a href="#">gihtub</a>
Python string similarity algorithm library		<a href="#">github</a>
Similar sentence judgment model based on Siamese bilstm model, providing training data set and test data set	100,000 training samples provided	<a href="#">github</a>

## Text Data Augmentation

Resource name (Name)	Description	Link
Chinese NLP Data Augmentation (EDA) Tool		<a href="#">github</a>
English NLP data enhancement tool		<a href="#">github</a>
One-click Chinese data enhancement tool		<a href="#">github</a>
The application and effect of data enhancement in machine translation and other nlp tasks		<a href="#">link</a>
NLP Data Augmentation Resource Collection		<a href="#">github</a>

# Common regular expressions

Resource name (Name)	Description	Link
Regular expression to extract email		It has been integrated into the python package <a href="#">cocoNLP</a> , welcome to try
Extract phone_number		It has been integrated into the python package <a href="#">cocoNLP</a> , welcome to try
Regular expression for extracting ID number	<pre>IDCards_pattern = r'^([1-9]\d{5}[12]\d{3}(0[1-9] 1[012])(0[1-9] 1[2][0-9] 3[01])\d{3}[0-9xX])'</pre> <pre>IDs = re.findall(IDCards_pattern, text, flags=0)</pre>	
IP address regular expression	<pre>(25[0-5]  2[0-4]\d  [0-1]\d{2}  [1-9]?\d).(25[0-5]  2[0-4]\d  [0-1]\d{2}  [1-9]?\d).(25[0-5]  2[0-4]\d  [0-1]\d {2}  [1-9]?\d).(25[0-5]  2[0-4]\d  [0-1]\d{2}  [1-9]?\d )</pre>	
Tencent QQ number regular expression	<pre>[1-9]([0-9]{5,11})</pre>	

Domestic  
fixed-line number  
regular expression

[0-9-()()]{7,18}

username regex [A-Za-z0-9\_-\u4e00-\u9fa5]+

Regular matching  
of domestic phone  
numbers (three  
major operators +  
virtual, etc.)

[github](#)

Regular  
Expression  
Tutorial

[github](#)

## text search

Resource name (Name)	Description	Link
Efficient Fuzzy Search Tool		<a href="#">github</a>
Large list/search engine of BERT models for various languages/tasks		<a href="#">link</a>
Deepmatch's deep matching model library for recommendation, advertising and search		<a href="#">github</a>

wwsearch is a full-text search engine developed by the enterprise WeChat background		<a href="#">github</a>
aili - the fastest in-memory index in the East		<a href="#">github</a>
Efficient string matching tool RapidFuzz	a fast string matching library for Python and C++, which is using the string similarity calculations from FuzzyWuzzy	<a href="#">github</a>

## reading comprehension

Resource name (Name)	Description	Link
Efficient Fuzzy Search Tool		<a href="#">github</a>
Large list/search engine of BERT models for various languages/tasks		<a href="#">link</a>
Deepmatch's deep matching model library for recommendation, advertising and search		<a href="#">github</a>
allennlp reading comprehension supports a variety of data and models		<a href="#">github</a>

## emotion analysis

Resource name (Name)	Description	Link
----------------------	-------------	------

aspect sentiment analysis package		<a href="#">github</a>
awesome-nlp-sentiment-analysis	Sentiment analysis, emotional cause identification, evaluation object and evaluation word extraction	<a href="#">github</a>
Sentiment analysis technology enables intelligent customer service to better understand human emotions		<a href="#">github</a>

# event extraction

Resource name (Name)	Description	Link
Chinese event extraction		<a href="#">github</a>
List of Literature Resources for NLP Event Extraction		<a href="#">github</a>
BERT event extraction implemented by PyTorch (ACE 2005 corpus)		<a href="#">github</a>
News Event Clue Extraction		<a href="#">github</a>

# machine translation

Resource name (Name)	Description	Link

no way dictionary	The command line version of Youdao Dictionary supports English-Chinese mutual search and online search	<a href="#">github</a>
NLLB	Language model NLLB that supports arbitrary inter-translation of 200+ languages	<a href="#">link</a>
Easy-Translate	Script to translate large text files locally, based on Facebook/Meta AI's M2M100 model and NLLB200 model, supports 200+ languages	<a href="#">github</a>

## digital conversion

Resource name (Name)	Description	Link
The best Chinese character number (Chinese number)-Arabic number conversion tool		<a href="#">github</a>
Quickly convert "Chinese numerals" and "Arabic numerals"		<a href="#">github</a>
Parse and convert natural language numeric strings to integers and floating point numbers		<a href="#">github</a>

## anaphora resolution

Resource name (Name)	Description	Link
Chinese reference to digestion data		<a href="#">github</a>

[baidu ink](#) code a0qq

## text clustering

Resource name (Name)	Description	Link
TextCluster short text clustering preprocessing module Short text cluster		<a href="#">github</a>

## Text Categorization

Resource name (Name)	Description	Link
NeuralNLP-NeuralClassifier Tencent open source deep learning text classification tool		<a href="#">github</a>

## knowledge reasoning

Resource name (Name)	Description	Link
GraphbrainAI is an open source software library and research tools designed to facilitate automatic meaning extraction and text understanding as well as knowledge exploration and inference		<a href="#">github</a>
(Harvard) free book on causal reasoning		<a href="#">pdf</a>

# Interpretable Natural Language Processing

Resource name (Name)	Description	Link
State-of-the-art interpreter library for textual machine learning models		<a href="#">github</a>

## text attack

Resource name (Name)	Description	Link
TextAttack natural language processing model adversarial attack framework		<a href="#">github</a>
OpenBackdoor: Text backdoor attack and defense toolkit	OpenBackdoor is developed based on Python and PyTorch, which can be used to reproduce, evaluate and develop related algorithms for text backdoor attack and defense	<a href="#">github</a>

## text visualization

Resource name (Name)	Description	Link
Scattertext text visualization (python)		<a href="#">github</a>



whatlies word vector interactive visualization		<a href="#">spacytools</a>
PySS3 machine visualization tool for SS3 text classifiers for explainable AI		<a href="#">github</a>
Render 3D images with Notepad		<a href="#">github</a>
attnvisGPT2, BERT and other transformer language models attention interactive visualization		<a href="#">github</a>
Textthero text data efficient processing package	Including preprocessing, keyword extraction, named entity recognition, vector space analysis, text visualization, etc.	<a href="#">github</a>

## text annotation tool

Resource name (Name)	Description	Link
Overview of NLP annotation platform		<a href="#">github</a>
brat rapid annotation tool sequence annotation tool		<a href="#">link</a>
Poplar web version natural language annotation tool		<a href="#">github</a>

LIDA is a lightweight interactive dialogue annotation tool

[github](#)

doccano is a web-based open source collaborative multilingual text annotation tool

[github](#)

Datasaurai online data labeling workflow management tool

[link](#)

## language detection

Resource name (Name)	Description	Link
langid	97 languages detected	<a href="https://github.com/saffsd/langid.py">https://github.com/saffsd/langid.py</a>
langdetect	language detection	<a href="https://code.google.com/archive/p/language-detection/">https://code.google.com/archive/p/language-detection/</a>

## comprehensive tool

Resource name (Name)	Description	Link
jieba		<a href="#">jieba</a>
hanlp		<a href="#">hanlp</a>

nlp4han	Chinese natural language processing tool set (sentence segmentation/word segmentation/part-of-speech tagging/chunking/syntax analysis/semantic analysis/NER/N-gram/HMM/pronoun resolution/sentiment analysis/spelling check	<a href="#">github</a>
Progress in Hate Speech Detection		<a href="#">link</a>
Bert application based on Pytorch	Including named entity recognition, sentiment analysis, text classification and text similarity, etc.	<a href="#">github</a>
nlp4han Chinese natural language processing toolset	Sentence segmentation/word segmentation/part-of-speech tagging/chunking/syntactic analysis/semantic analysis/NER/N-gram/HMM/pronoun resolution/sentiment analysis/spelling check	<a href="#">github</a>
Some basic models of natural language		<a href="#">github</a>
Template code for sequence tagging and text classification with BERT		<a href="#">github</a>
jieba_fast accelerated version of jieba		<a href="#">github</a>
Stanford NLP	Pure Python version of natural language processing package	<a href="#">link</a>

Python Spoken  
Natural Language  
Processing Toolset  
(English) [github](#)

PreNLP natural  
language  
preprocessing  
library [github](#)

Some papers and  
codes related to nlp Including topic model, word vector (Word  
Embedding), named entity recognition  
(NER), text classification (Text Classificatin),  
text generation (Text Generation), text  
similarity (Text Similarity) calculation, etc.,  
involving various nlp-related Algorithm,  
based on keras and tensorflow [github](#)

Python text  
mining/NLP practical  
example [github](#)

Forte's flexible and  
powerful natural  
language  
processing pipeline  
toolset [github](#)

stanza Stanford  
team NLP tools Can handle more than sixty languages [github](#)

Fancy-NLP is a text  
knowledge mining  
tool for building  
product portraits [github](#)

Comprehensive and  
easy Chinese NLP  
toolkit

[github](#)

Recurrence of  
vectorized recall  
pipelines commonly  
used in the industry  
based on DSSM

[github](#)

Textthero text data  
efficient processing  
package

Including preprocessing, keyword extraction,  
named entity recognition, vector space  
analysis, text visualization, etc.

[github](#)

nlpgnn graph neural  
network natural  
language  
processing toolbox

[github](#)

Macadam

Based on Tensorflow (Keras) and  
bert4keras, a natural language processing  
toolkit focusing on text classification,  
sequence labeling and relation extraction

[github](#)

LineFlow is an  
efficient NLP data  
loader for all deep  
learning frameworks

[github](#)

Arabica: Python text  
data exploratory  
analysis toolkit

[github](#)

Python stress  
testing tool:  
SMSBoom

[github](#)

# funny tool

Resource name (Name)	Description	Link
Wang Feng Lyric Generator		<a href="#">phunterlau/wangfeng-rnn</a>
Analysis of girlfriend's emotional fluctuations		<a href="#">github</a>
NLP is too difficult series		<a href="#">github</a>
Variable naming artifact		<a href="#">github link</a>
Image text removal, can be used for manga translation		<a href="#">github</a>
CoupletAI - couplet generation	Automatic couplet system based on CNN+Bi-LSTM+Attention	<a href="#">github</a>
Solving Complex Mathematical Equations Using Neural Network Symbolic Reasoning		<a href="#">github</a>

Question answering robot based on 14W song knowledge base	Functions include Lyrics Solitaire, Finding Songs with Known Lyrics, and Questions and Answers about the Triangular Relationship of Song Artists Lyrics	<a href="#">github</a>
COPE - Metric Poem Editor		<a href="#">github</a>
Paper2GUI	An AI desktop APP toolbox for ordinary people. It can be used immediately without installation. It already supports 18+ AI models, covering speech synthesis, video frame complementing, video super-resolution, target detection, image stylization, OCR recognition, etc.	<a href="#">github</a>
Politeness estimator (trained using Sina Weibo data)		<a href="#">github</a> <a href="#">paper</a>
Grass python (Python Chinese version) getting started guide	Chinese programming language	<a href="#">homepage</a> <a href="#">gitee</a>

## course report interview

Resource name (Name)	Description	Link
-------------------------	-------------	------

Natural Language  
Processing Report

[link](#)

Knowledge Graph  
Report

[link](#)

Data Mining Report

[link](#)

autonomous driving  
report

[link](#)

Machine translation  
report

[link](#)

blockchain report

[link](#)

robot report

[link](#)

Computer Graphics  
Report

[link](#)

3D printing report

[link](#)

Facial Recognition  
Report

[link](#)

Artificial Intelligence  
Chip Report

[link](#)

cs224n deep learning  
natural language  
processing course

pytorch  
implementation of the  
model in the [link](#)  
[courselink](#)



Natural Language  
Processing by  
Example Tutorial for  
Deep Learning  
Researchers

[github](#)

"Natural Language  
Processing" by Jacob  
Eisenstein

[github](#)

ML-NLP

Machine learning (Machine  
Learning), knowledge  
points and code  
implementation often  
tested in NLP interviews

[github](#)

NLP task example  
project code set

[github](#)

2019 NLP Highlights  
Review

[download](#)

nlp-recipes produced  
by Microsoft--best  
practices and  
examples of natural  
language processing

[github](#)

Natural Language  
Processing by  
Example Tutorial for  
Deep Learning  
Researchers

[github](#)

Transfer Learning in  
Natural Language  
Processing (NLP)

[youtube](#)

Machine Learning  
Systems book

[link](#) [github](#)

## Contest

Resource name (Name)	Description	Link
Review the TOP solutions of all NLP competitions		<a href="#">github</a>
2019 Baidu Triple Extraction Competition, "Scientific Space Team" source code (7th place)		<a href="#">github</a>

## Financial Natural Language Processing

Resource name (Name)	Description	Link
BDCI2019 Financial Negative Information Judgment		<a href="#">github</a>
Open source financial investment data extraction tool		<a href="#">github</a>
A large list of natural language processing research resources in the financial field		<a href="#">github</a>
Chatbots based on the financial-judicial domain (with the nature of small talk)		<a href="#">github</a>

Demonstration of small-scale financial knowledge  
graph construction process

[github](#)

## Medical Natural Language Processing

Resource name (Name)	Description	Link
Chinese medical NLP public resources arrangement		<a href="#">github</a>
spaCy Medical Text Mining and Information Extraction		<a href="#">github</a>
Building a Model for Medical Entity Recognition	Contains dictionaries and corpus annotations, based on python	<a href="#">github</a>
Question answering system based on knowledge graph in medical field		<a href="#">github</a> This repo refers to <a href="#">github</a>
Chinese medical dialogue data Chinese medical dialogue data set		<a href="#">github</a>
A Large-Scale Medical Dialogue Dataset	Contains 1.1 million medical consultations and 4 million doctor-patient dialogues	<a href="#">github</a>

Data related to  
COVID-19

New crown and other types of  
pneumonia Chinese medical  
dialogue dataset; open data  
sources of institutions such as  
Tsinghua University (COVID-19)

[github](#)

[github](#)

## Legal Natural Language Processing

Resource name (Name)	Description	Link
Blackstone's spaCy pipeline and NLP model for unstructured legal text		<a href="#">github</a>
List of Forensic Intelligence Literature Resources		<a href="#">github</a>
Chatbots based on the financial-judicial domain (with the nature of small talk)		<a href="#">github</a>
Crime Legal Terms and Classification Model	Contains 856 crime knowledge graphs, crime prediction based on 2.8 million crime training database, 13 types of question classification and legal information question and answer function based on 20W legal question and answer pairs	<a href="#">github</a>

text to image

Resource name (Name)	Description	Link
Dalle-mini	A mini version of DALL·E that generates pictures based on text prompts	<a href="#">github</a>

## other

Resource name (Name)	Description	Link
phone	China mobile phone attribution query	<a href="#">ls0f/phone</a>
phone	International mobile phone and telephone attribution inquiry	<a href="#">AfterShip/phone</a>
ngender	gender based on name	<a href="#">observers/ngender</a>
A summary of the differences between Chinese and English natural language processing NLP		<a href="#">link</a>
Technical documents PDF or PPT shared by Daniel in each major company		<a href="#">github</a>
comparxiv is used to compare the difference between two submitted versions on arXiv		<a href="#">pypi</a>
Meta-architecture of CHAMELEON deep learning news recommendation system		<a href="#">github</a>
Automatic Resume Screening System		<a href="#">github</a>

A variety of text readability  
evaluation indicators  
implemented by Python

[github](#)

## **Data Science ML Full Stack Roadmap**

<https://github.com/hemansnation/Data-Science-ML-Full-Stack-2022>

**Join the Data Science & ML Full Stack WhatsApp Group Community here:  
If the group is full, please join another one.**

<https://chat.whatsapp.com/B7Mdp6QTMJ0KZYGWrziT3Y>  
<https://chat.whatsapp.com/HWDSJU4KXrXJlcn5Npp3Gm>  
<https://chat.whatsapp.com/DmATV5uaVY7IKrTMHDiHnr>  
<https://chat.whatsapp.com/Blz2n8QYSgdKWfQbJZxHtJ>

## **Join Telegram for Data Science ML AI Resources:**

<https://t.me/+sREuRiFssMo4YWJl>

Join Community on LinkedIn:

<https://www.linkedin.com/groups/12540639/>

## **Connect with me on these platforms:**

LinkedIn: <https://www.linkedin.com/in/hemansnation/>

Twitter: <https://twitter.com/hemansnation>

GitHub: <https://github.com/hemansnation>

Instagram: <https://www.instagram.com/masterdexter.ai/>

## **Are you a professional?**

DM for One-on-One sessions for Python, Data Science, Machine Learning,  
and Data Engineering.

Here: <https://bit.ly/3U6zQvQ>

## **Python Notion Template**

<https://hemansnation.gumroad.com/l/god-level-python-with-himanshu-rachandani>