

# 《喜羊羊与灰太狼》知识图谱构建

本项目以动画《喜羊羊与灰太狼》为例，构建其中角色的人物关系知识图谱。项目任务一共分为四个部分：

- 1) 数据爬取与预处理
- 2) 命名实体识别
- 3) 实体关系抽取
- 4) 可视化与知识问答

对应项目代码分别在 [喜羊羊与灰太狼/数据爬取](#)，[喜羊羊与灰太狼/实体识别](#)，[喜羊羊与灰太狼/关系抽取](#)，[喜羊羊与灰太狼/可视化与知识问答](#) 目录中。

## 1 数据爬取与预处理

在构建知识图谱之前，需要收集项目需要的文本数据，从而进行实体（人物、场地）和关系的识别抽取。主要从《喜羊羊与灰太狼》维基百科网站和百度梗概文本获得信息，对应项目代码在 [喜羊羊与灰太狼/数据爬取](#) 目录下。

### 1.1 半结构化数据获取

我们选择维基百科（[喜羊羊与灰太狼 - 羊羊百科](#)，[青青草原的百科全书 - 灰机wiki \(huijiwiki.com\)](#)）中的半结构化数据文本作为爬取的目标。结构输出为包含两个实体一个关系的三元组格式，作为后面利用抽取技术结果的补充和比对。

#### 1.1.1 实体获取

由于每个角色的个人页面都是“<https://xyy.huijiwiki.com/wiki/角色名>”格式，我们选择在“角色目录”下爬取所有出场人物，并保存到 `entity.txt` 中，具体代码在 [喜羊羊与灰太狼/数据爬取/crawler\\_entity.py](#) 目录下。

```
session = HTMLSession()
response = session.get(url)
a_list = response.html.find('a')
cur_ents = []
for a in a_list:
    href_ent = a.attrs.get('href', '')
    if '/wiki/%' in href_ent:
        if('title' in a.attrs):
            cur_ents.append(a.attrs['title'])
for t in cur_ents:
    all_entities.append(t)
```

## 1.1.2 三元组获取

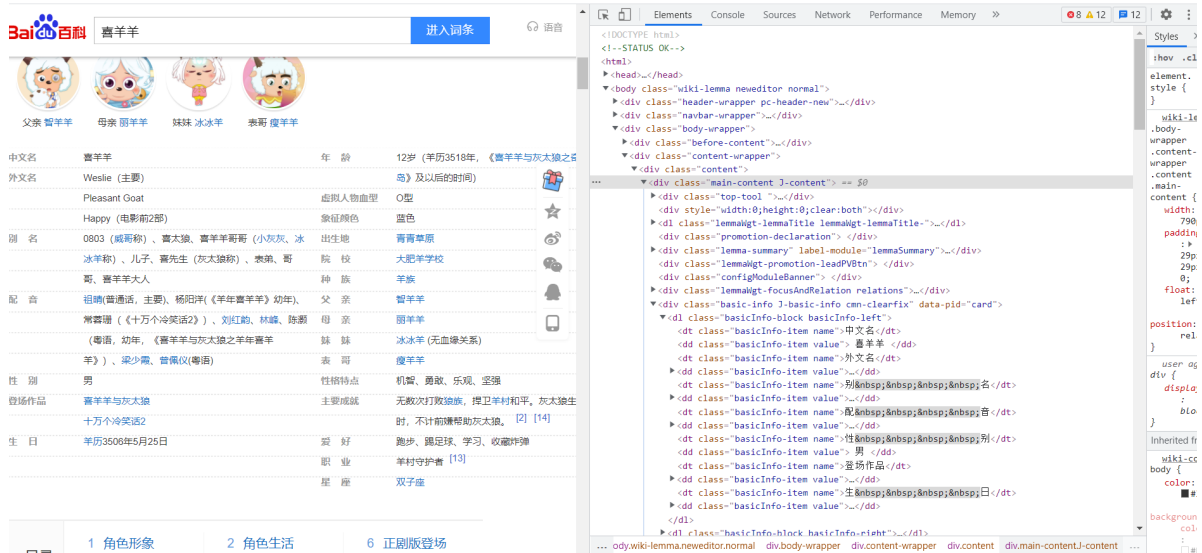
获取所有可以访问主页的结构后，对其半结构化数据进行分析，并从中抽取关系三元组，具体代码在 `喜羊羊与灰太狼/数据爬取/crawler_data.py` 目录下。

首先，检查wiki角色主页源代码，查看人物关系源代码格式如下：

接着，根据人物关系源代码格式，利用Beautiful Soup工具进行数据爬取，调用 `extract_entity(ent)` 函数用于实现，并将关系三元组爬取结果保存到 `data_raw.txt` 中。

```
def extract_entity(ent):
    url = root + ent
    response = requests.get(url=url,headers=headers).content
    soup = BeautifulSoup(response, 'html.parser')
    e1 = ent
    relations = []
    if soup.find('table'):
        tb = soup.find('table')
        trs = tb.findChildren('tr')
        for tr in trs:
            ths = tr.find_all('th',class_ = "infobox-label" )
            tds = tr.find_all('td',class_ = "infobox-data")
            for t in ths:
                t1 = t.text
                for t in tds:
                    t2 = t.text
                    t2 = t2.replace("\n", "")
                    r = (t1,t2)
                    relations.append(r)
    relations = [(e1, r[0], r[1]) for r in relations]
    return relations
```

另外，考虑到百度百科中也存在一些《喜羊羊与灰太狼》的人物关系信息，我们对其进行二次爬取以对三元组关系进行补充。其中，百度百科角色人物网址格式为“`https://baike.baidu.com/item/角色名`”，检查百度百科角色主页源代码，查看人物关系源代码格式如下：



调用 `baidu_rec()` 函数，对《喜羊羊与灰太狼》中的角色属性进行爬取，并将结果保存到 `角色名.txt` 文件中，具体代码在 `喜羊羊与灰太狼/数据爬取/bdbk_rec.py` 目录下。

```
def baidu_rec(root, seeds):
    while len(current_category) >= 1 :
        seed = current_category.pop(0)
        print('visiting', seed)
        url = root + seed
        session = HTMLSession()

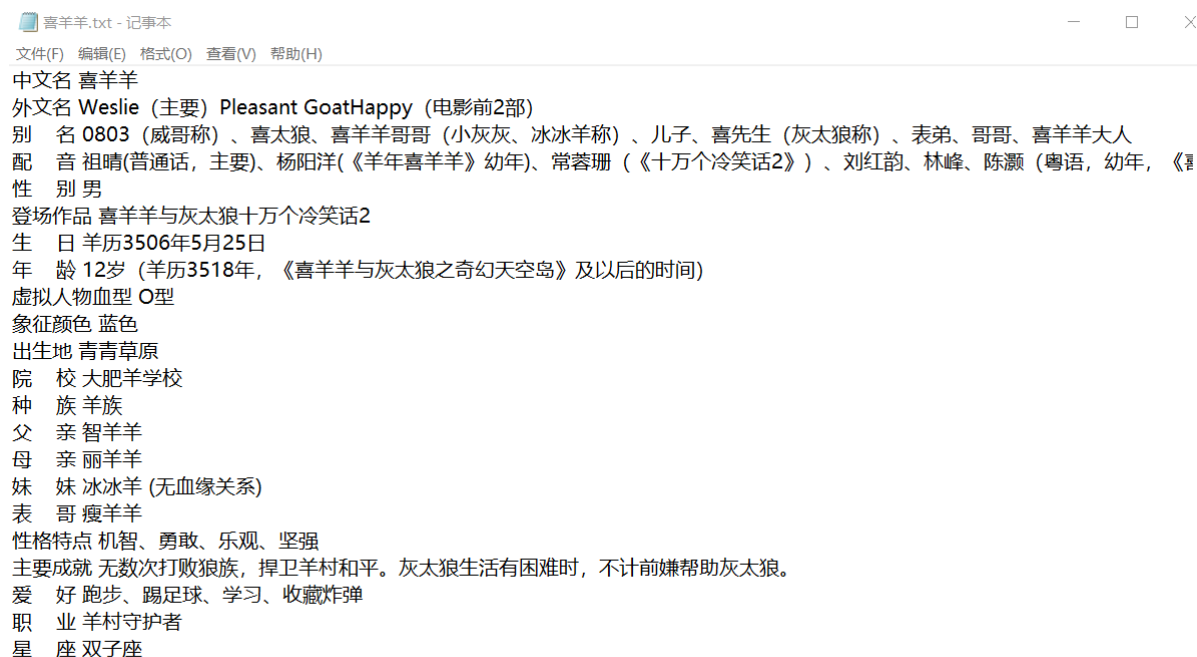
        response = session.get(url)

        html_code = response.content.decode(response.encoding)
        from lxml import etree
        treeObj = etree.HTML(html_code)
        target4 = treeObj.xpath("string(//dl[@class='basicInfo-block basicInfo-
left'])")
        print(len(target4))
        selector = Selector(text=html_code)
        title = ''.join(selector.xpath('//h1/text()').extract()).replace('/',
        '')
        names = selector.xpath('//dt[contains(@class,"basicInfo-item
name")]').extract()
        values = selector.xpath('//dd[contains(@class,"basicInfo-item
value")]').extract()

        lines = ''
        for i, name in enumerate(names) :
            # name
            temp =
            selector(text=name).xpath('//dt/text()//dt/a/text()').extract()
            name = ''.join(temp).replace('\n', '')
            # value
            temp =
            selector(text=values[i]).xpath('//dd/text()//dd/a/text()').extract()
            value = ''.join(temp).replace('\n', '')
            lines += name + '$$' + value + '\n'
            print(name, value)
        print('process file:' + str(title))
        output = open(title + '.txt', 'w', encoding='utf-8')
        output.write(lines)
```

```
output.close()
```

例如，喜羊羊在百度百科爬取关系结果如下：



喜羊羊.txt - 记事本

文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)

中文名 喜羊羊

外文名 **Weslie** (主要) **Pleasant GoatHappy** (电影前2部)

别名 0803 (威哥称)、喜太狼、喜羊羊哥哥 (小灰灰、冰冰羊称)、儿子、喜先生 (灰太狼称)、表弟、哥哥、喜羊羊大人

配音 祖晴(普通话, 主要)、杨阳洋(《羊年喜羊羊》幼年)、常蓉珊 (《十万个冷笑话2》)、刘红韵、林峰、陈颢 (粤语, 幼年, 《喜羊羊与灰太狼之牛气冲天》幼年)

性别 男

登场作品 喜羊羊与灰太狼十万个冷笑话2

生日 羊历3506年5月25日

年龄 12岁 (羊历3518年, 《喜羊羊与灰太狼之奇幻天空岛》及以后的时间)

虚拟人物血型 O型

象征颜色 蓝色

出生地 青青草原

院校 大肥羊学校

种族 羊族

父亲 智羊羊

母亲 丽羊羊

妹妹 冰冰羊 (无血缘关系)

表哥 瘦羊羊

性格特点 机智、勇敢、乐观、坚强

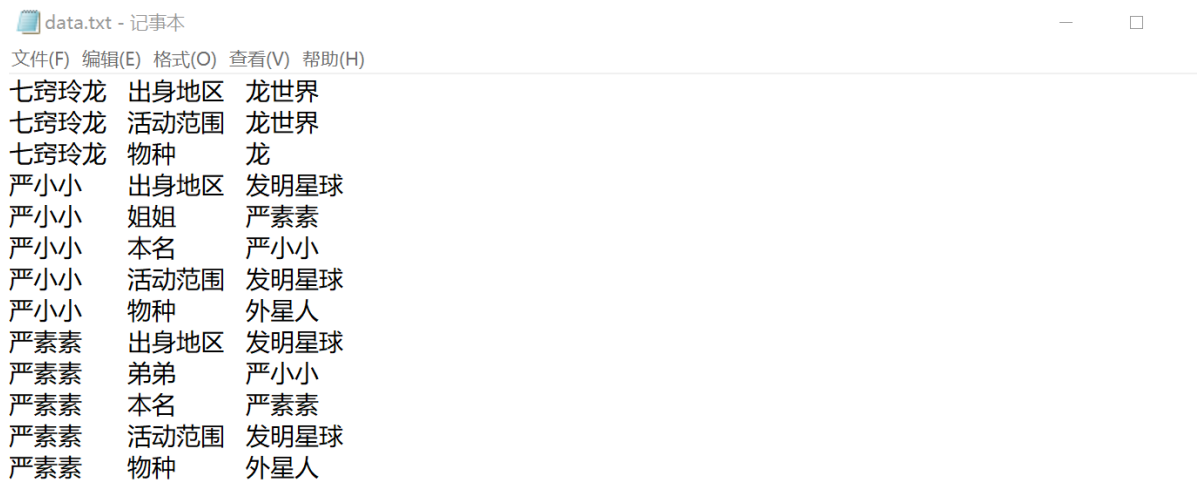
主要成就 无数次打败狼族, 捍卫羊村和平。灰太狼生活有困难时, 不计前嫌帮助灰太狼。

爱好 跑步、踢足球、学习、收藏炸弹

职业 羊村守护者

星座 双子座

另外，由于爬取数据中存在需要清洗的部分，需要对半结构化数据进行两部分清洗。一方面是，删掉不需要的例如配音演员、发行日期等信息；另一方面，处理不正确的符号，不正确的三元组格式，错别字等。将清洗结束后的最终结果保存在 `喜羊羊与灰太狼/数据爬取/data.txt`。



data.txt - 记事本

文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)

七窍玲龙 出身地区 龙世界

七窍玲龙 活动范围 龙世界

七窍玲龙 物种 龙

严小小 出身地区 发明星球

严小小 姐姐 严素素

严小小 本名 严小小

严小小 活动范围 发明星球

严小小 物种 外星人

严素素 出身地区 发明星球

严素素 弟弟 严小小

严素素 本名 严素素

严素素 活动范围 发明星球

严素素 物种 外星人

最后，把清洗好的三元组数据输出为csv格式，具体代码在 `喜羊羊与灰太狼/数据爬取/tocsv_data.py`。

```
with open(relations_final, 'w', encoding='utf-8') as f:
    for d in data:
        f.write('%s,%s,%s\n'%(d[0],d[1],d[2]))
```

最终关系三元组结果保存至 `喜羊羊与灰太狼/数据爬取/data.csv`，部分结果如下所示：

	A	B	C	D	E	F	G	H	I	J	K
1	七窍玲珑	龙世界	出身地区								
2	七窍玲珑	龙世界	活动范围								
3	七窍玲珑	龙	物种								
4	严小小	发明星球	出身地区								
5	严小小	严素素	姐姐								
6	严小小	严小小	本名								
7	严小小	发明星球	活动范围								
8	严小小	外星人	物种								
9	严素素	发明星球	出身地区								
10	严素素	严小小	弟弟								
11	严素素	严素素	本名								
12	严素素	发明星球	活动范围								
13	严素素	外星人	物种								
14	YY	懒羊羊	朋友								
15	YY	YY	本名								

## 1.2 文本数据获取

由于基于深度学习方法抽取人物关系，需要包含实体和关系的文本数据。我们选择百度中《喜羊羊与灰太狼》动漫每一集的剧情介绍 (<https://wenku.baidu.com/view/2d0edebdac51f01dc281e53a580216fc700a539d.html>) 作为文本，爬取方式与上文相同，并保存至 `喜羊羊与灰太狼/数据爬取/data.txt` 中。

由于爬取下来的文本数据存在一些错别字和不正确的符号，例如“灰人狼”，“喜洋洋”，所以先对文本数据进行简单清洗，最后得到

动画片400集数据量的训练数据，结果如下所示。

```

data.txt - 记事本
文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)
1. 《狼来了》(上)
相传，青青草原住着世界上最肥的羊，灰太狼带着祖先留下的典籍，携同妻子红太狼 来到青青草原，
2. 《狼来了》(下)
为了对抗灰太狼，慢羊羊决定重新启用实验室，众羊羊纷纷回忆起当年慢慢的发明把 实验室弄得一团
3. 《大小药丸》
灰太狼用变大丸与变小丸，几经周折才成功穿越羊村闸门，逮住懒羊羊。聪明的懒羊 羊用激将法令得
4. 《昏睡果》
灰太狼发现了一种昏睡果，便将其磨成汁，制出昏睡枪，灰太狼试验于红太狼，红太 狼睡着。灰太狼
5. 《变色狼》
为纪念慢羊羊村长功积，众羊合力造了一个蜡像，前往蜡像途中经狼家附近。灰太狼 像变色龙一样，

```

## 2 命名实体识别

为了更好地把数据中爬取的文本作为关系抽取的输入，我们利用命名实体识别技术对它进行处理，提取出包含两个实体以上的句子作为关系抽取的输入。通过比较 BI-LSTM 模型以及 BI-LSTM+CRF 模型对实体的识别效果，选择识别精确度相对最高的BI-LSTM+CRF模型作为本实验场景的最终实体识别模型。对应项目代码在 `喜羊羊与灰太狼/实体识别` 目录下。

## 2.1 训练数据标注

由于现有的NER训练数据都是以人名为主，不适合《喜羊羊与灰太狼》中的动物角色命名，所以需要手动进行数据标注。常见的标注方法有BIO、BIOE、BIOES等，这里我们选择BIOE编码方法。通过标注工具辅助，最终得到训练数据100条，将其输出保存在 喜羊羊与灰太狼/实体识别/entity/train.txt 目录下，部分结果如下：

 train.txt - 记事本

文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)

```
当 O
灰 B-Person
太 I-Person
狼 E-Person
的 O
身 O
影 O
再 O
次 O
```

## 2.2 基于BI-LSTM的命名实体识别

采用常见的 LSTM 对命名实体进行识别，具体代码为 喜羊羊与灰太狼/实体识别/entity/lstm.py。

首先，调用 read\_data() 和 build\_vocab() 函数，构造 word\_to\_id 和 tag\_to\_id 两个词典，将训练数据的文字序列和标注序列都转为数字数组，利用pytorch自带的词嵌入对文本序列进行预处理。

```
def read_data():
    sentences = []
    tags = []
    with open('data/train.txt', 'r', encoding='utf-8') as f:
        tmp_sentence = []
        tmp_tags = []
        for line in f:
            if line == '\n' and len(tmp_sentence) != 0:
                assert len(tmp_sentence) == len(tmp_tags)
                sentences.append(tmp_sentence)
                tags.append(tmp_tags)
                tmp_sentence = []
                tmp_tags = []
            elif line == '\n' or line == ' \n' or line == ' \n' :
                continue
            else:
                line = line.strip().split(' ')
                tmp_sentence.append(line[0])
                tmp_tags.append(line[1])
        if len(tmp_sentence) != 0:
            assert len(tmp_sentence) == len(tmp_tags)
            sentences.append(tmp_sentence)
            tags.append(tmp_tags)
    return sentences, tags

def build_vocab(sentences):
    global word_to_id
    for sentence in sentences: # 建立word到索引的映射
```

```

for word in sentence:
    if word not in word_to_id:
        word_to_id[word] = len(word_to_id)
return word_to_id

```

其次，通过LSTM层和分类器，输出tag词典大小的数组，得到的数据是某个标识符可能的最大概率，并利用交叉熵和梯度下降法优化，模型框架部分代码如下：

```

class Mymodel(nn.Module):
    def __init__(self, corpus_num, embedding_num, hidden_num, class_num, bi=True):
        super().__init__()
        self.embedding = nn.Embedding(corpus_num, embedding_num)
        self.lstm =
nn.LSTM(embedding_num, hidden_num, batch_first=True, bidirectional=bi)
        if bi :
            self.classifier = nn.Linear(hidden_num * 2, class_num)
        else:
            self.classifier = nn.Linear(hidden_num, class_num)
        self.cross_loss = nn.CrossEntropyLoss()

    def forward(self, batch_data, batch_tag=None):
        embedding = self.embedding(batch_data)
        out,_ = self.lstm(embedding)
        pre = self.classifier(out)
        self.pre = torch.argmax(pre, dim=-1).reshape(-1)
        #找到对应的最大的那个tag
        if batch_tag is not None:
            loss =
self.cross_loss(pre.reshape(-1, pre.shape[-1]), batch_tag.reshape(-1))
            #结果和已知做loss
        return loss

```

但是训练结果并不理想，会出现一些BIO逻辑上的错误，例如以I开头的实体，这显然不符合标注常理。

## 2.3 基于BI-LSTM+CRF的命名实体识别

由于LSTM模型实体识别效果并不理想，我们引入了概率论中的条件随机场CRF，对原LSTM模型进行改进。首先，把双向LSTM训练得到的结果作为发射矩阵，CRF作为转移矩阵，然后通过计算每次给定序列得分并找到最大得分路径的方式，最终构建BI-LSTM+CRF模型。具体代码在 `喜羊羊与灰太狼/实体识别/entity/lstm-crf.py` 目录下。

首先，根据CRF的定义，设计辅助函数 `log_sum_exp` 计算分数，通过调用 `log_sum_exp()` 函数进行实现：

```

def log_sum_exp(vec):
    max_score = vec[0, argmax(vec)]
    max_score_broadcast = max_score.view(1, -1).expand(1, vec.size()[1])
    return max_score + \
        torch.log(torch.sum(torch.exp(vec - max_score_broadcast)))

```

接着，构建BI-LSTM+CRF模型，部分代码如下：

```

def _forward_alg(self, feats):
    init_alphas = torch.full((1, self.tagset_size), -10000.)

```



```

init_alphas[0][self.tag_to_ix[START_TAG]] = 0.

forward_var = init_alphas

for feat in feats:
    alphas_t = []
    for next_tag in range(self.tagset_size):
        emit_score = feat[next_tag].view(
            1, -1).expand(1, self.tagset_size)
        trans_score = self.transitions[next_tag].view(1, -1)
        next_tag_var = forward_var + trans_score + emit_score
        alphas_t.append(log_sum_exp(next_tag_var).view(1))
    forward_var = torch.cat(alphas_t).view(1, -1)
terminal_var = forward_var + self.transitions[self.tag_to_ix[STOP_TAG]]
alpha = log_sum_exp(terminal_var)
return alpha

def _get_lstm_features(self, sentence):
    self.hidden = self.init_hidden()
    embeds = self.word_embeds(sentence).view(len(sentence), 1, -1)
    lstm_out, self.hidden = self.lstm(embeds, self.hidden)
    lstm_out = lstm_out.view(len(sentence), self.hidden_dim)
    lstm_feats = self.hidden2tag(lstm_out)
    return lstm_feats

def _score_sentence(self, feats, tags):

    score = torch.zeros(1)
    tags = torch.cat([torch.tensor([self.tag_to_ix[START_TAG]],
dtype=torch.long), tags])
    for i, feat in enumerate(feats):
        score = score + \
            self.transitions[tags[i + 1], tags[i]] + feat[tags[i + 1]]
    score = score + self.transitions[self.tag_to_ix[STOP_TAG], tags[-1]]
    return score

def _viterbi_decode(self, feats):
    backpointers = []

    init_vvars = torch.full((1, self.tagset_size), -10000.)
    init_vvars[0][self.tag_to_ix[START_TAG]] = 0

    forward_var = init_vvars
    for feat in feats:
        bptrs_t = []
        viterbivars_t = []

        for next_tag in range(self.tagset_size):
            next_tag_var = forward_var + self.transitions[next_tag]
            best_tag_id = argmax(next_tag_var)
            bptrs_t.append(best_tag_id)
            viterbivars_t.append(next_tag_var[0][best_tag_id].view(1))
        forward_var = (torch.cat(viterbivars_t) + feat).view(1, -1)
        backpointers.append(bptrs_t)

    terminal_var = forward_var + self.transitions[self.tag_to_ix[STOP_TAG]]
    best_tag_id = argmax(terminal_var)
    path_score = terminal_var[0][best_tag_id]

```



```

best_path = [best_tag_id]
for bptrs_t in reversed(backpointers):
    best_tag_id = bptrs_t[best_tag_id]
    best_path.append(best_tag_id)
start = best_path.pop()
assert start == self.tag_to_ix[START_TAG]
best_path.reverse()
return path_score, best_path

def neg_log_likelihood(self, sentence, tags):
    feats = self._get_lstm_features(sentence)
    forward_score = self._forward_alg(feats)
    gold_score = self._score_sentence(feats, tags)
    return forward_score - gold_score

def forward(self, sentence):
    lstm_feats = self._get_lstm_features(sentence)

    score, tag_seq = self._viterbi_decode(lstm_feats)
    return score, tag_seq

```

最后，根据训练数据 `train.txt` 对模型进行训练，并基于维特比解码预测全部数据，得到最终实体识别结果。对数据再次进行清洗后，将文本 BIOE 编码结果保存至 `喜羊羊与灰太狼\实体识别\entity\data\result.txt` 中。

```

result.txt - 记事本
文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)
山羊大叔因小乖乖太捣蛋将其托给众羊照顾，众羊被小乖乖的恶作剧吓躲到地下室，灰太狼乘机捉走小乖乖，问地
B-Person I-Person I-Person E-Person O B-Person I-Person E-Person O O O O O O O O O O O O O O B-P
沸羊羊与懒羊羊自诩为有义气，鄙视喜羊羊没义气，然而当他们被灰太狼抓住时，却是喜羊羊不计前嫌，冒着危险
B-Person I-Person E-Person O B-Person I-Person E-Person O O O O O O O O B-Person I-Person E-Persc
每当慢羊羊上山采药时，总要带一名助手帮忙。懒羊羊用尽千方百计后，终于得到助手一职。但笨手笨脚的懒羊羊
O O B-Person I-Person E-Person O O O O O O O O O O O O O O B-Person I-Person E-Person O O O

```

## 3 实体关系抽取

为抽取出文本中的实体关系，首先，基于开源 OpenNRE 工具包 (<https://github.com/thunlp/OpenNRE>) 的数据输入格式，对实体识别后的 BIOE 编码结果进行格式化处理。其次，调用 `opennre` 库，分别训练 `wiki80_cnn_softmax`、`wiki80_bert_softmax`、`wiki80_bertentity_softmax`、`tacred_bert_softmax` 以及 `tacred_bertentity_softmax`，5 种关系抽取的预训练模型。接着，调用哈工大，BERT-wwm，中文 bert，在 20w 中文人物关系数据，训练 `chinese-BERT-wwm` 模型 (<https://github.com/taorui-plus/OpenNRE>)。最后，将《喜羊羊与灰太狼》格式化处理后的数据，输入 6 种预训练模型进行推理，输出关系抽取结果。对应项目代码在 `喜羊羊与灰太狼\关系抽取` 目录下。

### 3.1 输入数据处理

由于 OpenNRE 数据对文本以及对应实体标注的输入格式如下：

```

`{'text': 'He was the son of Máel Dúin mac Máele Fithrich, and grandson of the high king Áed Uaridnach (died 612).', 'h': {'pos': (18, 46)}, 't': {'pos': (78, 91)}}`

```

因此，需要将实体识别后的数据格式进行转换，运行 `喜羊羊与灰太狼\关系抽取/predata.py`，得到处理后的输入数据 `new_data`。

```
def read_data():
    sentences = []
    tags = []
    i=-1
    with open('result.txt', 'r', encoding='gbk') as f:
        line = f.readline()          # 调用文件的 readline()方法
        while line:
            i=i+1

            if len(line)==0 or line=='\n':
                line = f.readline()
                continue
            if i%3==0 or i==0:
                sentences.append(line)
                line = f.readline()
            else:
                tags.append(line)
                line = f.readline()

    return sentences, tags
sentences, tags=read_data()
new_data=[]

for t in range(len(tags)):
    tags[t]=tags[t].replace('-Person','')
    ans=tags[t].replace(' ','')
    #print(len(ans),len(sentences[t]))
    blist=[]
    elist=[]
    for a in re.finditer('B', ans):
        blist.append(a.span())
    for a in re.finditer('E', ans):
        elist.append(a.span())
    if (len(blist)!=len(elist)):
        for k in range(1,len(blist)-1):
            if (blist[k][0]-blist[k-1][1])==0:
                del blist[k-1]
        if (len(blist)!=len(elist)):
            print(blist)
            print(elist)
            print(sentences[t])

    bb=[]
    for kk in range(len(blist)-1):
        aa=[]
        aa.append(blist[kk][0])
        aa.append(elist[kk][0])
        bb.append(aa)
    for ss in range(len(bb)):
        for tt in range(ss+1,len(bb)):
            p1=bb[ss][0]
            q1=bb[ss][1]
```

```
p2=bb[tt][0]
q2=bb[tt][1]
posstr1= '('+str(p1)+','+str(q1)+')'
posstr2= '('+str(p2)+','+str(q2)+')'
new_data.append({'text':sentences[t].replace('\n',''), 'h':
{'pos':eval(posstr1)}, 't': {'pos':eval(posstr2)}})
```

第一次运行后，结果输出如下，发现存在处理异常的数据。

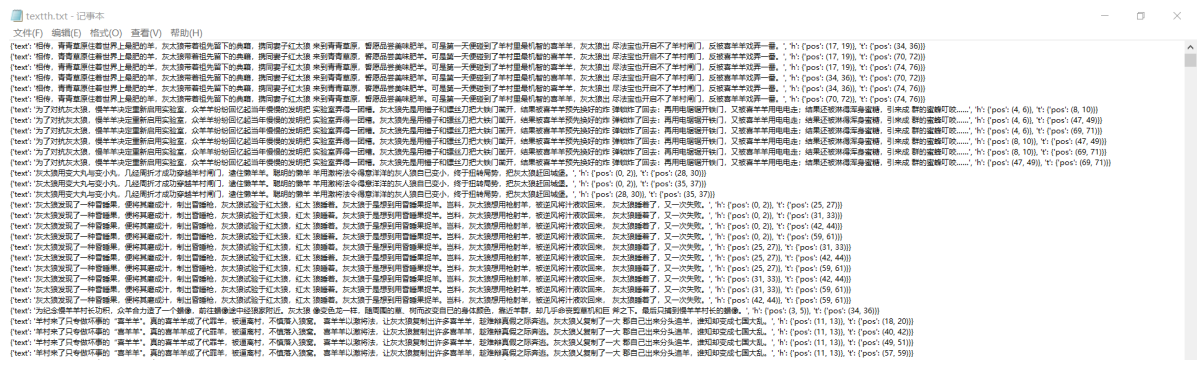
```
[(12, 13), (21, 22), (46, 47), (61, 62), (85, 86)]
[(14, 15), (23, 24), (48, 49), (63, 64), (77, 78), (87, 88)]
```

新年到，众羊要交换礼物，懒羊羊把口水垫给了暖羊羊，可是他又后悔了，要想尽办法拿回来，结果被灰太狼捉住了，还当成礼物送给了狼表哥，但很快又后悔了，于是与表哥 争抢起来，让喜羊羊他们有机可乘，把他们弹飞了。

根据输出信息，对实体识别结果 result.txt 再次进行检查，发现实体识别将“表哥”标注为“I-Person E-Person”，导致进行上述处理时无法匹配到“表哥”实体的开始。因为此处“表哥”代指“狼表哥”，且已经出现在前面语句中，因此手动将“表哥”实体标注更改为“o o”。然后再次运行上述代码，并将结果写入 喜羊羊与灰太狼\关系抽取/textth.txt 文件中。

```
f = open("textth.txt", "w", errors='ignore')
for line in new_data:
    f.write(str(line))
    f.write('\n')
f.close()
```

打开 textth.txt 文件，数据处理结果如下：



打印其中一条处理结果：

```
{'text': '相传，青青草原住着世界上最肥的羊，灰太狼带着祖先留下的典籍，携同妻子红太狼 来到青青草原，誓愿品尝美味肥羊。可是第一天便碰到了羊村里最机智的喜羊羊，灰太狼出 尽法宝也开启不了羊村闸门，反被喜羊羊戏弄一番。',
'h': {'pos': (17, 19)},
't': {'pos': (70, 72)}}
```

### 3.2 基于opennre的关系抽取

处理完关系抽取模型的输入数据后，接下来通过opennre工具对 wiki80\_cnn\_softmax、wiki80\_bert\_softmax、wiki80\_bertentity\_softmax、taced\_bert\_softmax 以及 taced\_bertentity\_softmax，这5种关系抽取模型进行预训练，然后对实验数据进行推理，得到最后的关系抽取结果。

首先，在windows10上安装opennre工具。打开cmd，运行以下指令，安装相关库以及包：

```
git clone https://github.com/thunlp/OpenNRE.git --depth 1
pip install -r requirements.txt
python setup.py install
```

安装完成后，尝试运行测试代码：

```
import opennre
model = opennre.get_model('wiki80_cnn_softmax').cuda()
```

代码报错，显示无法找到 `pretrain/glove/glove.6B.50d_word2id.json` 文件。这是因为，在代码自动运行 `download_glove.sh` 文件时，`wget glove.6B.50d_word2id.json` 文件链接显示“连接失败”，需要手动查找对应文件的有效下载地址：

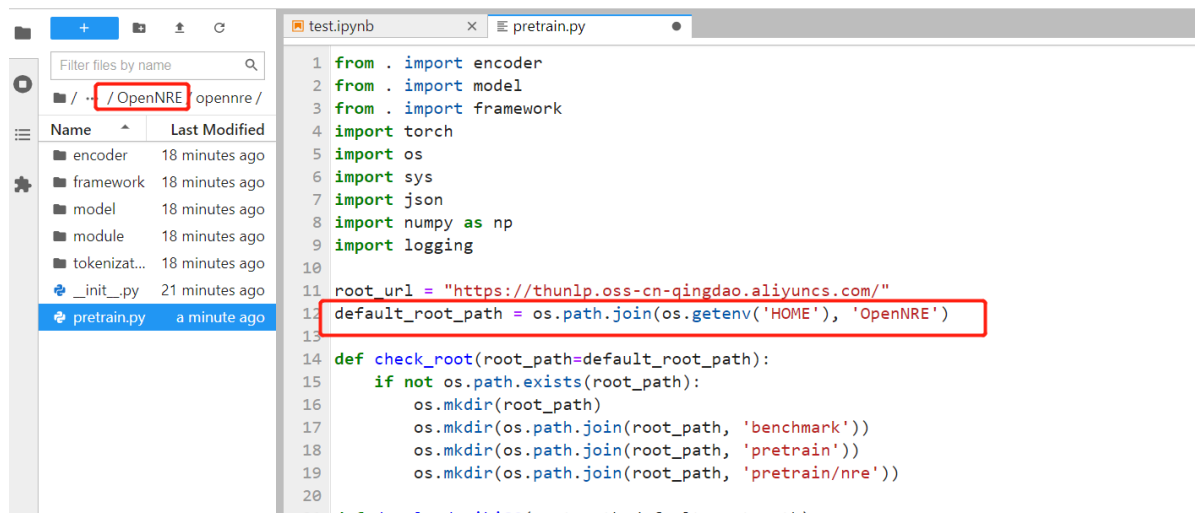
```
https://thunlp.oss-cn-qingdao.aliyuncs.com/opennre/pretrain/glove/glove.6B.50d_word2id.json
```

保存 `glove.6B.50d_word2id.json` 文件至 `opennre\pretrain/glove` 目录下，同理保存 `glove.6B.50d_mat.npy` 文件。再次运行，发现仍然找不到对应文件。debug进入 `./opennre/pretrain.py` 文件，发现引用文件路径需要配置环境变量。

将 `pretrain.py` 文件中：

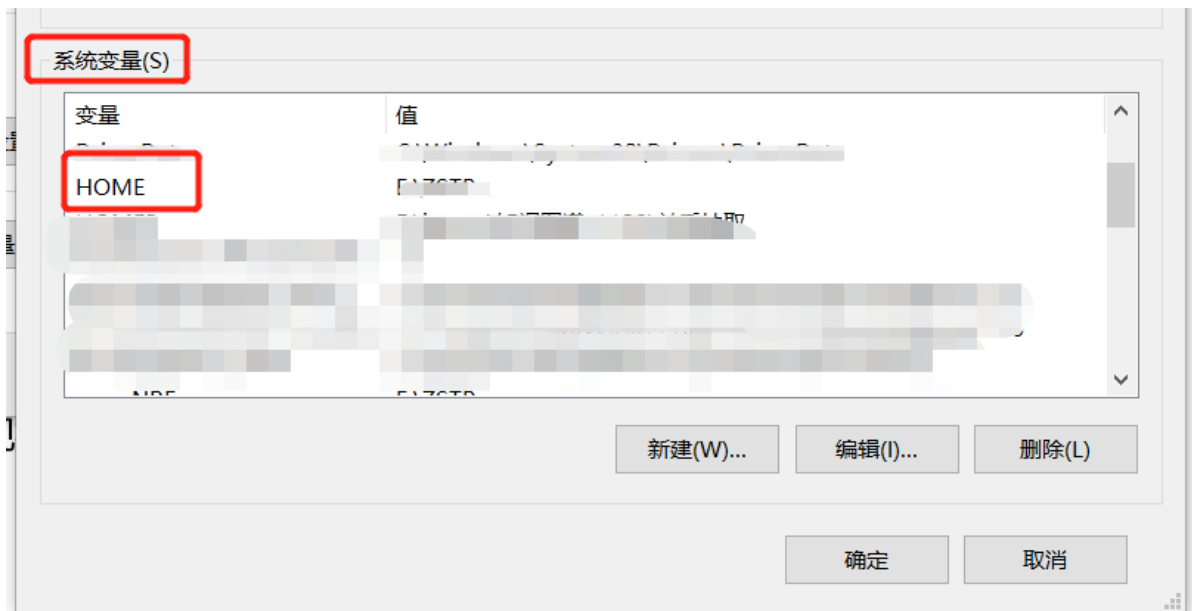
```
default_root_path = os.path.join(os.getenv('HOME'), '.opennre')
```

根据项目名称更改为：



```
default_root_path = os.path.join(os.getenv('HOME'), 'OpenNRE')
```

同时，将 `'HOME'` 以及对项目所在路径添加进环境变量。



再次运行测试代码，发现缺少 `benchmark/wiki80/wiki80_rel2id.json`、`benchmark/wiki80/wiki80_train.txt` 以及 `pretrain/nre/wiki80_cnn_softmax.pth.tar` 文件，同理找到下载链接，对缺失文件进行下载。

```
https://thunlp.oss-cn-qingdao.aliyuncs.com/opennre/benchmark/wiki80/wiki80_rel2id.json
https://thunlp.oss-cn-qingdao.aliyuncs.com/opennre/benchmark/wiki80/wiki80_train.txt
https://thunlp.oss-cn-qingdao.aliyuncs.com/opennre/pretrain/nre/wiki80_cnn_softmax.pth.tar
```

再次运行测试代码，代码运行成功，输出相应信息：

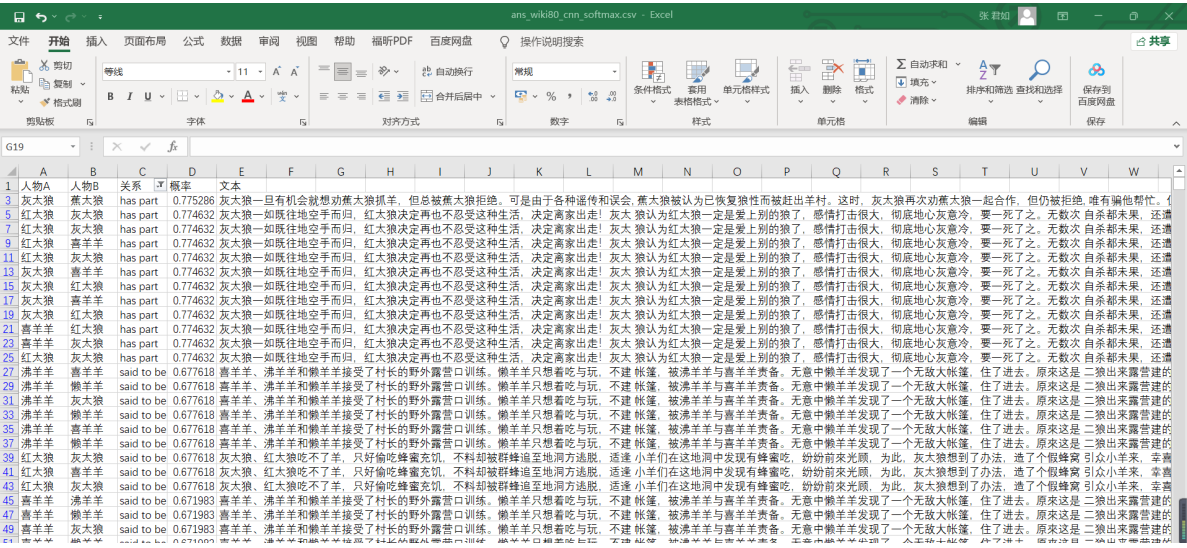
```
2021-11-29 23:59:49,889 - root - INFO - Initializing word embedding with word2vec.
```

opennre库测试成功后，首先，将刚才调用成功的 `wiki80_cnn_softmax` 模型，代入 `喜羊羊与灰太狼\关系抽取\OpenNRE/relation_train.py`：

```
import csv
fw=open('ans_wiki80_cnn_softmax.csv','w')
csv_writer = csv.writer(fw)
csv_writer.writerow(["人物A","人物B","关系","概率","文本"])
with open('textth.txt', 'r', encoding='gb18030',errors='ignore') as f :
    for line in f :
        if line[0].isdigit() :
            continue
        elif len(line) != 0 :
            ss=eval(line)
            #print(ss['text'][ss['h']['pos'][0]:ss['h']['pos'][1]+1], ss['text']
[ss['t']['pos'][0]:ss['t']['pos'][1]+1])
            result = model.infer(ss)
            relation=result[0].encode('gb18030').decode('utf-8')
            #print(relation)
            a=ss['text'][ss['h']['pos'][0]:ss['h']['pos'][1]+1]
            b=ss['text'][ss['t']['pos'][0]:ss['t']['pos'][1]+1]
            if a==b:
                continue
```

```
else:
    csv_writer.writerow([a, b, relation, result[1], ss['text']])
print("training")
fw.close()
```

输出 ans\_wiki80\_cnn\_softmax.csv 文件:



从模型运行结果上看，发现关系抽取结果并不理想，例如将“沸羊羊”与“喜羊羊”识别为同一个人物：

ID	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W
1	人物A	人物B	关系	概率	文本																		
19	灰太狼	红太狼	has part	0.774632	灰太狼—如既往地空手而归。红太狼决定再也不想受这种生活，决定离家出走！	灰太狼认为红太狼一定是爱上别的狼了，感情打击很大，彻底地心灰意冷。																	
21	喜羊羊	灰太狼	has part	0.774632	灰太狼—如既往地空手而归。红太狼决定再也不想受这种生活，决定离家出走！	灰太狼认为红太狼一定是爱上别的狼了，感情打击很大，彻底地心灰意冷。																	
23	喜羊羊	灰太狼	has part	0.774632	灰太狼—如既往地空手而归。红太狼决定再也不想受这种生活，决定离家出走！	灰太狼认为红太狼一定是爱上别的狼了，感情打击很大，彻底地心灰意冷。																	
25	红太狼	灰太狼	has part	0.774632	灰太狼—如既往地空手而归。红太狼决定再也不想受这种生活，决定离家出走！	灰太狼认为红太狼一定是爱上别的狼了，感情打击很大，彻底地心灰意冷。																	
26	沸羊羊	喜羊羊	said to be	0.677618	喜羊羊、沸羊羊和懒羊羊接受了村长的野外露营地训练。懒羊羊只想看吃与玩。不建帐篷，被沸羊羊与喜羊羊责备。无意中懒羊羊发现了一个无敌大帐篷，住了进去。原来是二狼出来露营建的。																		
29	沸羊羊	懒羊羊	said to be	0.677618	喜羊羊、沸羊羊和懒羊羊接受了村长的野外露营地训练。懒羊羊只想看吃与玩。不建帐篷，被沸羊羊与喜羊羊责备。无意中懒羊羊发现了一个无敌大帐篷，住了进去。原来是二狼出来露营建的。																		
31	沸羊羊	灰太狼	said to be	0.677618	喜羊羊、沸羊羊和懒羊羊接受了村长的野外露营地训练。懒羊羊只想看吃与玩。不建帐篷，被沸羊羊与喜羊羊责备。无意中懒羊羊发现了一个无敌大帐篷，住了进去。原来是二狼出来露营建的。																		
35	沸羊羊	喜羊羊	said to be	0.677618	喜羊羊、沸羊羊和懒羊羊接受了村长的野外露营地训练。懒羊羊只想看吃与玩。不建帐篷，被沸羊羊与喜羊羊责备。无意中懒羊羊发现了一个无敌大帐篷，住了进去。原来是二狼出来露营建的。																		
37	沸羊羊	懒羊羊	said to be	0.677618	喜羊羊、沸羊羊和懒羊羊接受了村长的野外露营地训练。懒羊羊只想看吃与玩。不建帐篷，被沸羊羊与喜羊羊责备。无意中懒羊羊发现了一个无敌大帐篷，住了进去。原来是二狼出来露营建的。																		
39	沸羊羊	灰太狼	said to be	0.677618	喜羊羊、沸羊羊和懒羊羊接受了村长的野外露营地训练。懒羊羊只想看吃与玩。不建帐篷，被沸羊羊与喜羊羊责备。无意中懒羊羊发现了一个无敌大帐篷，住了进去。原来是二狼出来露营建的。																		
41	红太狼	喜羊羊	said to be	0.677618	喜羊羊、沸羊羊和懒羊羊接受了村长的野外露营地训练。懒羊羊只想看吃与玩。不建帐篷，被沸羊羊与喜羊羊责备。无意中懒羊羊发现了一个无敌大帐篷，住了进去。原来是二狼出来露营建的。																		
43	红太狼	灰太狼	said to be	0.671983	喜羊羊、沸羊羊和懒羊羊接受了村长的野外露营地训练。懒羊羊只想看吃与玩。不建帐篷，被沸羊羊与喜羊羊责备。无意中懒羊羊发现了一个无敌大帐篷，住了进去。原来是二狼出来露营建的。																		
45	喜羊羊	沸羊羊	said to be	0.671983	喜羊羊、沸羊羊和懒羊羊接受了村长的野外露营地训练。懒羊羊只想看吃与玩。不建帐篷，被沸羊羊与喜羊羊责备。无意中懒羊羊发现了一个无敌大帐篷，住了进去。原来是二狼出来露营建的。																		

出现这样的结果，考虑是由于 wiki80\_cnn\_softmax 模型是基于wiki80数据集，使用 CNN 编码器对数据集进行训练，而wiki80数据集主要是针对英文语句的词向量间关系进行抽取的，因此并不太适用于中文文本，不能很好地解析中文文本的语义信息。另外，由于CNN模型很难根据训练的结果去针对性的调整具体的特征，大多只能根据词向量之间的距离对实体关系进行推断，因此导致对于本实验场景的关系抽取结果并不理想。

接着，尝试调用 wiki80\_bert\_softmax 模型:

```
model = openre.get_model('wiki80_bert_softmax').cuda()
```

报错缺少相关配置，下载 bert 模型，并配置至 pretrain/bert-base-uncased。

- <https://thunlp.oss-cn-qingdao.aliyuncs.com/openre/pretrain/bert-base-uncased/config.json>
- [https://thunlp.oss-cn-qingdao.aliyuncs.com/openre/pretrain/bert-base-uncased/pytorch\\_model.bin](https://thunlp.oss-cn-qingdao.aliyuncs.com/openre/pretrain/bert-base-uncased/pytorch_model.bin)
- <https://thunlp.oss-cn-qingdao.aliyuncs.com/openre/pretrain/bert-base-uncased/vocab.txt>

再次调用 wiki80\_bert\_softmax 模型，对实验数据进行关系抽取:

```
import openre
model = openre.get_model('wiki80_bert_softmax').cuda()
import csv
fw=open('ans_wiki80_bert_softmax.csv','w')
```



```

csv_writer = csv.writer(fw)
csv_writer.writerow(["人物A", "人物B", "关系", "概率", "文本"])
with open('textth.txt', 'r', encoding='gb18030', errors='ignore') as f :
    for line in f :
        if line[0].isdigit() :
            continue
        elif len(line) != 0 :
            ss=eval(line)
            #print(ss['text'][ss['h']]['pos'][0]:ss['h']]['pos'][1]+1], ss['text']
            [ss['t']]['pos'][0]:ss['t']]['pos'][1]+1])
            result = model.infer(ss)
            relation=result[0].encode('gb18030').decode('utf-8')
            #print(relation)
            a=ss['text'][ss['h']]['pos'][0]:ss['h']]['pos'][1]+1]
            b=ss['text'][ss['t']]['pos'][0]:ss['t']]['pos'][1]+1]
            if a==b:
                continue
            else:
                csv_writer.writerow([a, b, relation,result[1],ss['text']])
            print("training")
fw.close()

```

输出 ans\_wiki80\_bert\_softmax.csv 文件:

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
1	人物A	人物B	关系	概率	文本															
3	巨蛙	灰太狼	followed by	0.992761	巨蛙村民被灰太狼勒令建祭坛。在小羊们的帮助下很快就建好了。但是灰太狼却要求他们加建。这时没有石料了。巨蛙和小羊们该怎么办呢？															
5	懒羊羊	灰太狼	followed by	0.991305	灰太狼控制着机器战警进入羊村并捉住了懒羊羊。喜羊羊把机器战警引到木桥上。掉到河里。救回懒羊羊。灰太狼再次将机器战警改良。改成用木头做。做到不少羊。喜羊羊															
7	巨蛙	灰太狼	followed by	0.99092	巨蛙村民被灰太狼勒令建祭坛。在小羊们的帮助下很快就建好了。但是灰太狼却要求他们加建。这时没有石料了。巨蛙和小羊们该怎么办呢？															
9	灰太狼	喜羊羊	followed by	0.990881	为了修葺教科书。村长去偷拍灰太狼照片而被捉。众小羊救出村长。但村长却按照死板的村规给众小羊嘉奖。众小羊不乐。村长再次拍照被捉。众小羊再救村却因争功导致															
11	喜羊羊	灰太狼	followed by	0.990867	灰太狼控制着机器战警进入羊村并捉住了懒羊羊。喜羊羊把机器战警引到木桥上。掉到河里。救回懒羊羊。灰太狼再次将机器战警改良。改成用木头做。做到不少羊。喜															
13	懒羊羊	灰太狼	followed by	0.990648	灰太狼控制着机器战警进入羊村并捉住了懒羊羊。喜羊羊把机器战警引到木桥上。掉到河里。救回懒羊羊。灰太狼再次将机器战警改良。改成用木头做。做到不少羊。喜															
15	懒羊羊	喜羊羊	followed by	0.990006	灰太狼控制着机器战警进入羊村并捉住了懒羊羊。喜羊羊把机器战警引到木桥上。掉到河里。救回懒羊羊。灰太狼再次将机器战警改良。改成用木头做。做到不少羊。喜															
17	灰太狼	懒羊羊	followed by	0.989732	懒羊羊想吃河中鲜草。伙伴们搭桥帮他过河。之后他竟把桥推倒。独占鲜草。不料那是灰太狼的陷阱。懒羊羊无人救援。束手就擒。喜羊羊称要除掉自私的懒羊羊。利用															
19	喜羊羊	懒羊羊	followed by	0.989707	懒羊羊想吃河中鲜草。伙伴们搭桥帮他过河。之后他竟把桥推倒。独占鲜草。不料那是灰太狼的陷阱。懒羊羊无人救援。束手就擒。喜羊羊称要除掉自私的懒羊羊。利用															
21	懒羊羊	喜羊羊	followed by	0.989562	一次无意的事故中。村长误以为是懒羊羊所为。其实是喜羊羊所作。任凭喜羊羊怎样百般讨好懒羊羊。就是得不到懒羊羊的原谅。喜羊羊感到烦恼。懒羊羊不听大家的劝告															
23	沸羊羊	喜羊羊	followed by	0.989488	沸羊羊和喜羊羊在很多方面都差不多。美羊羊都处处偏袒喜羊羊。见沸羊羊很苦恼。喜羊羊想出假装灰太狼帮助沸羊上演英雄救美的办法。当沸羊羊却遇上了真正义															
25	红太狼	灰太狼	followed by	0.989475	春天来了。燕子等鸟类开始筑巢。灰太狼夫妇故意捉弄燕子。却被其夺走红太狼戒指上的钻石。灰太狼为了取回钻石。不分昼夜意图害上燕子巢。发现此事的众羊。出于义															
27	懒羊羊	喜羊羊	followed by	0.989414	一次无意的事故中。村长误以为是懒羊羊所为。其实是喜羊羊所作。任凭喜羊羊怎样百般讨好懒羊羊。就是得不到懒羊羊的原谅。喜羊羊感到烦恼。懒羊羊不听大家的劝告															
29	懒羊羊	灰太狼	followed by	0.989364	一次无意的事故中。村长误以为是懒羊羊所为。其实是喜羊羊所作。任凭喜羊羊怎样百般讨好懒羊羊。就是得不到懒羊羊的原谅。喜羊羊感到烦恼。懒羊羊不听大家的劝告															
31	懒羊羊	喜羊羊	followed by	0.989364	一次无意的事故中。村长误以为是懒羊羊所为。其实是喜羊羊所作。任凭喜羊羊怎样百般讨好懒羊羊。就是得不到懒羊羊的原谅。喜羊羊感到烦恼。懒羊羊不听大家的劝告															
33	喜羊羊	懒羊羊	followed by	0.989307	喜羊羊因为出色的表现。而招致沸羊羊和懒羊羊的嫉妒。两只小羊合谋趁喜羊羊他们郊游的机会。利用灰太狼来打击喜羊羊。不料弄巧成拙。两羊反被灰所捉。幸好喜羊															
35	灰太狼	喜羊羊	followed by	0.989288	灰太狼控制着机器战警进入羊村并捉住了懒羊羊。喜羊羊把机器战警引到木桥上。掉到河里。救回懒羊羊。灰太狼再次将机器战警改良。改成用木头做。做到不少羊。喜															
37	懒羊羊	喜羊羊	followed by	0.989283	一次无意的事故中。村长误以为是懒羊羊所为。其实是喜羊羊所作。任凭喜羊羊怎样百般讨好懒羊羊。就是得不到懒羊羊的原谅。喜羊羊感到烦恼。懒羊羊不听大家的劝告															
39	懒羊羊	喜羊羊	followed by	0.989256	喜羊羊因为出色的表现。而招致沸羊羊和懒羊羊的嫉妒。两只小羊合谋趁喜羊羊他们郊游的机会。利用灰太狼来打击喜羊羊。不料弄巧成拙。两羊反被灰所捉。幸好喜羊															
41	灰太狼	喜羊羊	followed by	0.989196	一次无意的事故中。村长误以为是懒羊羊所为。其实是喜羊羊所作。任凭喜羊羊怎样百般讨好懒羊羊。就是得不到懒羊羊的原谅。喜羊羊感到烦恼。懒羊羊不听大家的劝告															
43	懒羊羊	懒羊羊	followed by	0.989196	一次无意的事故中。村长误以为是懒羊羊所为。其实是喜羊羊所作。任凭喜羊羊怎样百般讨好懒羊羊。就是得不到懒羊羊的原谅。喜羊羊感到烦恼。懒羊羊不听大家的劝告															
45	懒羊羊	灰太狼	followed by	0.989149	村长研制出一条小拉锁。能拉开任何物品。然而这个宝物却被灰太狼在懒羊羊手上骗了去。灰太狼利用这拉锁拉开墙壁。企图进入羊村。可是拉锁太小。灰太狼被卡在墙中															
47	喜羊羊	懒羊羊	followed by	0.989086	一次无意的事故中。村长误以为是懒羊羊所为。其实是喜羊羊所作。任凭喜羊羊怎样百般讨好懒羊羊。就是得不到懒羊羊的原谅。喜羊羊感到烦恼。懒羊羊不听大家的劝告															
49	美羊羊	灰太狼	followed by	0.989051	草原持续干旱。美羊羊却因为洗澡浪费大量清水。灰太狼趁机。做成一间流动浴室。将美羊羊擒获。却因为烧掉美羊羊时引发了火灾。喜羊羊等以德报怨。用村杖存下的															
51	灰太狼	懒羊羊	followed by	0.989016	懒羊羊想吃河中鲜草。伙伴们搭桥帮他过河。之后他竟把桥推倒。独占鲜草。不料那是灰太狼的陷阱。懒羊羊无人救援。束手就擒。喜羊羊称要除掉自私的懒羊羊。利用															
53	懒羊羊	灰太狼	followed by	0.988943	一次无意的事故中。村长误以为是懒羊羊所为。其实是喜羊羊所作。任凭喜羊羊怎样百般讨好懒羊羊。就是得不到懒羊羊的原谅。喜羊羊感到烦恼。懒羊羊不听大家的劝告															
55	懒羊羊	喜羊羊	followed by	0.988912	一次无意的事故中。村长误以为是懒羊羊所为。其实是喜羊羊所作。任凭喜羊羊怎样百般讨好懒羊羊。就是得不到懒羊羊的原谅。喜羊羊感到烦恼。懒羊羊不听大家的劝告															

从抽取结果上看，由于 wiki80\_bert\_softmax 模型是基于 wiki80 数据集，使用 BERT 编码器对数据集进行训练，相对于 wiki80\_cnn\_softmax 模型有所提升，但是对于人物关系的抽取并不是很理想，考虑原因仍是和训练数据集的数据特点有关。

同理，调用 taced\_bert\_softmax 模型，并下载相关配置文件，输出

ans\_taced\_bert\_softmax.csv 文件结果如下：



D3	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W
1	人物A	人物B	关系	概率	文本																		
3	灰太狼	红太狼	NA	0.995048	相传,青青草原住着世界上最肥的羊,灰太狼带着祖先留下的典籍,携同妻子红太狼来到青青草原,誓愿品尝美味肥羊。可是第一天便碰到了羊村里最机警的喜羊羊,灰太狼出尽法宝也开不了																		
4	灰太狼	喜羊羊	NA	0.997629	相传,青青草原住着世界上最肥的羊,灰太狼带着祖先留下的典籍,携同妻子红太狼来到青青草原,誓愿品尝美味肥羊。可是第一天便碰到了羊村里最机警的喜羊羊,灰太狼出尽法宝也开不了																		
7	红太狼	喜羊羊	NA	0.997381	相传,青青草原住着世界上最肥的羊,灰太狼带着祖先留下的典籍,携同妻子红太狼来到青青草原,誓愿品尝美味肥羊。可是第一天便碰到了羊村里最机警的喜羊羊,灰太狼出尽法宝也开不了																		
9	红太狼	喜羊羊	NA	0.997615	相传,青青草原住着世界上最肥的羊,灰太狼带着祖先留下的典籍,携同妻子红太狼来到青青草原,誓愿品尝美味肥羊。可是第一天便碰到了羊村里最机警的喜羊羊,灰太狼出尽法宝也开不了																		
11	喜羊羊	灰太狼	NA	0.997726	相传,青青草原住着世界上最肥的羊,灰太狼带着祖先留下的典籍,携同妻子红太狼来到青青草原,誓愿品尝美味肥羊。可是第一天便碰到了羊村里最机警的喜羊羊,灰太狼出尽法宝也开不了																		
13	灰太狼	慢羊羊	NA	0.996354	为了对抗灰太狼,慢羊羊决定重新启用实验室。众羊羊幼幼回忆起当年慢慢的发明把实验室弄得一团糟,灰太狼先用锤子和螺丝刀把大门凿开,结果被喜羊羊预先换好的炸弹锁住了回去,再																		
15	灰太狼	喜羊羊	NA	0.998031	为了对抗灰太狼,慢羊羊决定重新启用实验室。众羊羊幼幼回忆起当年慢慢的发明把实验室弄得一团糟,灰太狼先用锤子和螺丝刀把大门凿开,结果被喜羊羊预先换好的炸弹锁住了回去,再																		
17	慢羊羊	灰太狼	NA	0.995808	为了对抗灰太狼,慢羊羊决定重新启用实验室。众羊羊幼幼回忆起当年慢慢的发明把实验室弄得一团糟,灰太狼先用锤子和螺丝刀把大门凿开,结果被喜羊羊预先换好的炸弹锁住了回去,再																		
19	慢羊羊	喜羊羊	NA	0.995939	为了对抗灰太狼,慢羊羊决定重新启用实验室。众羊羊幼幼回忆起当年慢慢的发明把实验室弄得一团糟,灰太狼先用锤子和螺丝刀把大门凿开,结果被喜羊羊预先换好的炸弹锁住了回去,再																		
21	灰太狼	喜羊羊	NA	0.997971	为了对抗灰太狼,慢羊羊决定重新启用实验室。众羊羊幼幼回忆起当年慢慢的发明把实验室弄得一团糟,灰太狼先用锤子和螺丝刀把大门凿开,结果被喜羊羊预先换好的炸弹锁住了回去,再																		
23	灰太狼	懒羊羊	NA	0.999093	灰太狼用变大头与变小丸,几经周折才成功穿越羊村大门,逮住懒羊羊。聪明的懒羊羊用激将法令得意洋洋的灰太狼自己变小,终于扭转局势,把灰太狼赶回城堡。																		
25	灰太狼	懒羊羊	NA	0.999066	灰太狼用变大头与变小丸,几经周折才成功穿越羊村大门,逮住懒羊羊。聪明的懒羊羊用激将法令得意洋洋的灰太狼自己变小,终于扭转局势,把灰太狼赶回城堡。																		
27	懒羊羊	喜羊羊	NA	0.999172	灰太狼用变大头与变小丸,几经周折才成功穿越羊村大门,逮住懒羊羊。聪明的懒羊羊用激将法令得意洋洋的灰太狼自己变小,终于扭转局势,把灰太狼赶回城堡。																		
29	灰太狼	红太狼	NA	0.844406	灰太狼发现了一种昏睡果,便将其磨成汁,制出昏睡枪,灰太狼试验于红太狼,红太狼睡着。灰太狼于是想到用昏睡果捉羊。岂料,灰太狼想用枪射羊,被逆风将汁液吹回来,灰太狼睡着了,又																		
31	灰太狼	红太狼	NA	0.793865	灰太狼发现了一种昏睡果,便将其磨成汁,制出昏睡枪,灰太狼试验于红太狼,红太狼睡着。灰太狼于是想到用昏睡果捉羊。岂料,灰太狼想用枪射羊,被逆风将汁液吹回来,灰太狼睡着了,又																		
33	红太狼	灰太狼	NA	0.982343	灰太狼发现了一种昏睡果,便将其磨成汁,制出昏睡枪,灰太狼试验于红太狼,红太狼睡着。灰太狼于是想到用昏睡果捉羊。岂料,灰太狼想用枪射羊,被逆风将汁液吹回来,灰太狼睡着了,又																		
35	红太狼	灰太狼	NA	0.987733	灰太狼发现了一种昏睡果,便将其磨成汁,制出昏睡枪,灰太狼试验于红太狼,红太狼睡着。灰太狼于是想到用昏睡果捉羊。岂料,灰太狼想用枪射羊,被逆风将汁液吹回来,灰太狼睡着了,又																		
37	慢羊羊	灰太狼	NA	0.989527	为纪念慢羊羊村长功绩,众羊羊合力造了一个雕像,前往雕像途经经家附近,灰太狼像变色龙一样,周围的草、树而改变自己的身体颜色,靠近羊群,却几乎命丧铁犁和巨斧之下。最后只捕																		
39	喜羊羊	灰太狼	NA	0.991016	羊村来了只专做坏事的“喜羊羊”。真的喜羊羊成了代罪羊,被逼离村,不情愿入狼窝。喜羊羊以激将法,让灰太狼复制出许多喜羊羊,趁难辨真假之际奔逃。灰太狼复制了一大帮自己出来分头																		
41	喜羊羊	灰太狼	NA	0.990423	羊村来了只专做坏事的“喜羊羊”。真的喜羊羊成了代罪羊,被逼离村,不情愿入狼窝。喜羊羊以激将法,让灰太狼复制出许多喜羊羊,趁难辨真假之际奔逃。灰太狼复制了一大帮自己出来分头																		
43	喜羊羊	灰太狼	NA	0.991908	羊村来了只专做坏事的“喜羊羊”。真的喜羊羊成了代罪羊,被逼离村,不情愿入狼窝。喜羊羊以激将法,让灰太狼复制出许多喜羊羊,趁难辨真假之际奔逃。灰太狼复制了一大帮自己出来分头																		
45	灰太狼	喜羊羊	NA	0.999137	羊村来了只专做坏事的“喜羊羊”。真的喜羊羊成了代罪羊,被逼离村,不情愿入狼窝。喜羊羊以激将法,让灰太狼复制出许多喜羊羊,趁难辨真假之际奔逃。灰太狼复制了一大帮自己出来分头																		
47	沸羊羊	喜羊羊	NA	0.998528	沸羊羊钓鱼时上了潜水水底的灰太狼,因之前多次谎报捉狼,这次人家不信他,等发现沸羊羊不见,他已被快下狼锅了。喜羊羊以一群饿羊引开心的灰太狼捉羊,当狼大口咬下肥羊时,去																		
49	沸羊羊	灰太狼	NA	0.994787	沸羊羊钓鱼时上了潜水水底的灰太狼,因之前多次谎报捉狼,这次人家不信他,等发现沸羊羊不见,他已被快下狼锅了。喜羊羊以一群饿羊引开心的灰太狼捉羊,当狼大口咬下肥羊时,去																		

从抽取结果上看, `tacred_bert_softmax` 模型是基于TACRED数据集,使用BERT编码器对数据集进行训练。对于本实验场景,关系抽取的输出结果大多是‘NA’。考虑到这是因为TACRED数据集主要涉及的是英语新闻通讯社和网络文本,所以无法对中文的人物关系以及动画场景做出识别。

同理,调用 `tacred_bertentity_softmax` 模型以及 `wiki80_bertentity_softmax` 模型,并下载相关配置文件,输出 `ans_tacred_bert_softmax.csv` 以及 `ans_wiki80_bertentity_softmax.csv` 文件结果分别如下:

D127	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W
1	人物A	人物B	关系	概率	文本																		
3	灰太狼	红太狼	NA	0.997707	相传,青青草原住着世界上最肥的羊,灰太狼带着祖先留下的典籍,携同妻子红太狼来到青青草原,誓愿品尝美味肥羊。可是第一天便碰到了羊村里最机警的喜羊羊,灰太狼出尽法宝也开不了																		
5	灰太狼	喜羊羊	NA	0.996171	相传,青青草原住着世界上最肥的羊,灰太狼带着祖先留下的典籍,携同妻子红太狼来到青青草原,誓愿品尝美味肥羊。可是第一天便碰到了羊村里最机警的喜羊羊,灰太狼出尽法宝也开不了																		
7	红太狼	喜羊羊	NA	0.996265	相传,青青草原住着世界上最肥的羊,灰太狼带着祖先留下的典籍,携同妻子红太狼来到青青草原,誓愿品尝美味肥羊。可是第一天便碰到了羊村里最机警的喜羊羊,灰太狼出尽法宝也开不了																		
9	红太狼	喜羊羊	NA	0.996205	相传,青青草原住着世界上最肥的羊,灰太狼带着祖先留下的典籍,携同妻子红太狼来到青青草原,誓愿品尝美味肥羊。可是第一天便碰到了羊村里最机警的喜羊羊,灰太狼出尽法宝也开不了																		
11	喜羊羊	灰太狼	NA	0.999208	相传,青青草原住着世界上最肥的羊,灰太狼带着祖先留下的典籍,携同妻子红太狼来到青青草原,誓愿品尝美味肥羊。可是第一天便碰到了羊村里最机警的喜羊羊,灰太狼出尽法宝也开不了																		
13	灰太狼	慢羊羊	NA	0.997987	为了对抗灰太狼,慢羊羊决定重新启用实验室。众羊羊幼幼回忆起当年慢慢的发明把实验室弄得一团糟,灰太狼先用锤子和螺丝刀把大门凿开,结果被喜羊羊预先换好的炸弹锁住了回去,再																		
15	灰太狼	喜羊羊	NA	0.996871	为了对抗灰太狼,慢羊羊决定重新启用实验室。众羊羊幼幼回忆起当年慢慢的发明把实验室弄得一团糟,灰太狼先用锤子和螺丝刀把大门凿开,结果被喜羊羊预先换好的炸弹锁住了回去,再																		
17	慢羊羊	喜羊羊	NA	0.998870	为了对抗灰太狼,慢羊羊决定重新启用实验室。众羊羊幼幼回忆起当年慢慢的发明把实验室弄得一团糟,灰太狼先用锤子和螺丝刀把大门凿开,结果被喜羊羊预先换好的炸弹锁住了回去,再																		
21	灰太狼	喜羊羊	NA	0.998980	为了对抗灰太狼,慢羊羊决定重新启用实验室。众羊羊幼幼回忆起当年慢慢的发明把实验室弄得一团糟,灰太狼先用锤子和螺丝刀把大门凿开,结果被喜羊羊预先换好的炸弹锁住了回去,再																		
23	灰太狼	懒羊羊	NA	0.996205	灰太狼用变大头与变小丸,几经周折才成功穿越羊村大门,逮住懒羊羊。聪明的懒羊羊用激将法令得意洋洋的灰太狼自己变小,终于扭转局势,把灰太狼赶回城堡。																		
25	灰太狼	懒羊羊	NA	0.996205	灰太狼用变大头与变小丸,几经周折才成功穿越羊村大门,逮住懒羊羊。聪明的懒羊羊用激将法令得意洋洋的灰太狼自己变小,终于扭转局势,把灰太狼赶回城堡。																		
27	懒羊羊	喜羊羊	NA	0.996205	灰太狼用变大头与变小丸,几经周折才成功穿越羊村大门,逮住懒羊羊。聪明的懒羊羊用激将法令得意洋洋的灰太狼自己变小,终于扭转局势,把灰太狼赶回城堡。																		
29	灰太狼	红太狼	NA	0.996487	灰太狼发现了一种昏睡果,便将其磨成汁,制出昏睡枪,灰太狼试验于红太狼,红太狼睡着。灰太狼于是想到用昏睡果捉羊。岂料,灰太狼想用枪射羊,被逆风将汁液吹回来,灰太狼睡着了,又																		
31	灰太狼	红太狼	NA	0.997214	灰太狼发现了一种昏睡果,便将其磨成汁,制出昏睡枪,灰太狼试验于红太狼,红太狼睡着。灰太狼于是想到用昏睡果捉羊。岂料,灰太狼想用枪射羊,被逆风将汁液吹回来,灰太狼睡着了,又																		
33	红太狼	灰太狼	NA	0.996515	灰太狼发现了一种昏睡果,便将其磨成汁,制出昏睡枪,灰太狼试验于红太狼,红太狼睡着。灰太狼于是想到用昏睡果捉羊。岂料,灰太狼想用枪射羊,被逆风将汁液吹回来,灰太狼睡着了,又																		
35	红太狼	灰太狼	NA	0.996517	灰太狼发现了一种昏睡果,便将其磨成汁,制出昏睡枪,灰太狼试验于红太狼,红太狼睡着。灰太狼于是想到用昏睡果捉羊。岂料,灰太狼想用枪射羊,被逆风将汁液吹回来,灰太狼睡着了,又																		
37	慢羊羊	灰太狼	NA	0.997723	为纪念慢羊羊村长功绩,众羊羊合力造了一个雕像,前往雕像途经经家附近,灰太狼像变色龙一样,周围的草、树而改变自己的身体颜色,靠近羊群,却几乎命丧铁犁和巨斧之下。最后只捕																		
39	喜羊羊	灰太狼	NA	0.998621	羊村来了只专做坏事的“喜羊羊”。真的喜羊羊成了代罪羊,被逼离村,不情愿入狼窝。喜羊羊以激将法,让灰太狼复制出许多喜羊羊,趁难辨真假之际奔逃。灰太狼复制了一大帮自己出来分头																		
41	喜羊羊	灰太狼	NA	0.996949	羊村来了只专做坏事的“喜羊羊”。真的喜羊羊成了代罪羊,被逼离村,不情愿入狼窝。喜羊羊以激将法,让灰太狼复制出许多喜羊羊,趁难辨真假之际奔逃。灰太狼复制了一大帮自己出来分头																		
43	喜羊羊	灰太狼	NA	0.998040	羊村来了只专做坏事的“喜羊羊”。真的喜羊羊成了代罪羊,被逼离村,不情愿入狼窝。喜羊羊以激将法,让灰太狼复制出许多喜羊羊,趁难辨真假之际奔逃。灰太狼复制了一大帮自己出来分头																		
45	灰太狼	喜羊羊	NA	0.999017	羊村来了只专做坏事的“喜羊羊”。真的喜羊羊成了代罪羊,被逼离村,不情愿入狼窝。喜羊羊以激将法,让灰太狼复制出许多喜羊羊,趁难辨真假之际奔逃。灰太狼复制了一大帮自己出来分头																		
47	沸羊羊	喜羊羊	NA	0.998296	沸羊羊钓鱼时上了潜水水底的灰太狼,因之前多次谎报捉狼,这次人家不信他,等发现沸羊羊不见,他已被快下狼锅了。喜羊羊以一群饿羊引开心的灰太狼捉羊,当狼大口咬下肥羊时,去																		
49	沸羊羊	灰太狼	NA	0.999491	沸羊羊钓鱼时上了潜水水底的灰太狼,因之前多次谎报捉狼,这次人家不信他,等发现沸羊羊不见,他已被快下狼锅了。喜羊羊以一群饿羊引开心的灰太狼捉羊,当狼大口咬下肥羊时,去																		
51	沸羊羊	喜羊羊	NA	0.998653	沸羊羊钓鱼时上了潜水水底的灰太狼,因之前多次谎报捉狼,这次人家不信他,等发现沸羊羊不见,他已被快下狼锅了。喜羊羊以一群饿羊引开心的灰太狼捉羊,当狼大口咬下肥羊时,去																		

C29	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W
1	人物A	人物B	关系	概率	文本																		
3	慢羊羊	懒羊羊	said to be the same as	0.880134	慢羊羊为帮助村里最没用的羊而发明了“无敌装置”,并送给了懒羊羊。懒羊羊凭装置打败灰太狼而沾沾自喜。灰太狼设法换掉了装置后向懒羊羊挑战,懒羊羊应战																		
5	灰太狼	喜羊羊	followed by	0.87897	灰太狼决定要抓只羊招待将到狼婆表母。灰太狼用红太狼乔装成慢羊羊造成被其捕获的假象,诱使其其它小羊与之交换,慢羊羊及时出现戳穿灰太狼的阴谋。灰太狼																		
7	灰太狼	喜羊羊	followed by	0.875244	灰太狼决定要抓只羊招待将到狼婆表母。灰太狼用红太狼乔装成慢羊羊造成被其捕获的假象,诱使其其它小羊与之交换,慢羊羊及时出现戳穿灰太狼的阴谋。灰太狼																		
9	灰太狼	喜羊羊	said to be the same as	0.863179	懒羊羊正声称固土被占,前来羊村寻求帮助。众羊羊以为当然,但在他每晚啼啼的歌声干扰下,只得帮忙。可到沼泽一看,占地的竟然是灰太狼,灰太狼说:“无计可施了,喜																		
11	灰太狼	喜羊羊	followed by	0.859981	灰太狼决定要抓只羊招待将到狼婆表母。灰太狼用红太狼乔装成慢羊羊造成被其捕获的假象,诱使其其它小羊与之交换,慢羊羊及时出现戳穿灰太狼的阴谋。灰太狼																		
13	灰太狼	喜羊羊	followed by	0.859981	灰太狼决定要抓只羊招待将到狼婆表母。灰太狼用红太狼乔装成慢羊羊造成被其捕获的假象,诱使其其它小羊与之交换,慢羊羊及时出现戳穿灰太狼的阴谋。灰太狼																		
15	灰太狼	狼婆表母	followed by	0.859981	灰太狼决定要抓只羊招待将到狼婆表母。灰太狼用红太狼乔装成慢羊羊造成被其捕获的假象,诱使其其它小羊与之交换,慢羊羊及时出现戳穿灰太狼的阴谋。灰太狼																		
17	灰太狼	喜羊羊	followed by	0.859981	灰太狼决定要抓只羊招待将到狼婆表母。灰太狼用红太狼乔装成慢羊羊造成被其捕获的假象,诱使其其它小羊与之交换,慢羊羊及时出现戳穿灰太狼的阴谋。灰太狼																		
19	灰太狼	狼婆表母	followed by	0.859981	灰太狼决定要抓只羊招待将到狼婆表母。灰太狼用红太狼乔装成慢羊羊造成被其捕获的假象,诱使其其它小羊与之交换,慢羊羊及时出现戳穿灰太狼的阴谋。灰太狼																		
21	灰太狼	喜羊羊	followed by	0.854822	灰太狼决定要抓只羊招待将到狼婆表母。灰太狼用红太狼乔装成慢羊羊造成被其捕获的假象,诱使其其它小羊与之交换,慢羊羊及时出现戳穿灰太狼的阴谋。灰太狼																		
23	灰太狼	喜羊羊	said to be the same as	0.846675	懒羊羊正声称固土被占,前来羊村寻求帮助。众羊羊以为当然,但在他每晚啼啼的歌声干扰下,只得帮忙。可到沼泽一看,占地的竟然是灰太狼,灰太狼说:“无计可施了,喜																		
25	灰太狼	喜羊羊	said to be the same as	0.839932	懒羊羊正声称固土被占,前来羊村寻求帮助。众羊羊以为当然,但在他每晚啼啼的歌声干扰下,只得帮忙。可到沼泽一看,占地的竟然是灰太狼,灰太狼说:“无计可施了,喜																		
27	灰太狼	喜羊羊	said to be the same as	0.836887	草原上出现了一群乌大军,他们到处旅行,前进路上的障碍尽数破坏。狼婆表母的灰太狼狼胆失败,便想方设法引开狼婆表母向前逃。狼婆表的突然出令																		
29	灰太狼	喜羊羊	said to be the same as	0.831215	在一次露露中,懒羊羊在羊村里地要地要进入负责轻快的烹调工作的女生组,但进入女生组后,才发现做女生的要求太多,因而被赶出女生宿舍,还被灰太狼捉去。																		
31	黑太帅	喜羊羊	followed by	0.829718	灰太狼夫妇投书了黑太帅,接受了黑太帅的命令,带领了羊村的蜜蜂去围捕喜羊羊他们,就在这么危急关头,喜羊羊他们却发现了黑太帅的鬼魂,那他们真的																		
33	灰太狼	喜羊羊	followed by	0.829704	为了修葺教科书,村长去偷偷灰太狼的羊皮而捉。众羊羊救出村长,但村长却按照灰太狼的计策给村长小羊喜羊羊,众羊羊不乐,村长再次拍狼被捉,众羊羊再次救村却因																		
35	灰太狼	喜羊羊	instance of	0.824201	灰太狼利用装置作为诱饵诱捕他人,却机灵的羊们选择了,为了救出装置,众羊羊决定向黑太帅假投降,让喜羊羊与黑太帅在蓝蓝中合作,正当计划顺利进行时																		
37	灰太狼	喜羊羊	followed by	0.820389	灰太狼决定要抓只羊招待将到狼婆表母。灰太狼用红太狼乔装成慢羊羊造成被其捕获的假象,诱使其其它小羊与之交换,慢羊羊及时出现戳穿灰太狼的阴谋。灰太狼																		
39	灰太狼	喜羊羊	followed by	0.820389	灰太狼决定要抓只羊招待将到狼婆表母。灰太狼用红太狼乔装成慢羊羊造成被其捕获的假象,诱使其其它小羊与之交换,慢羊羊及时出现戳穿灰太狼的阴谋。灰太狼																		
41	灰太狼	喜羊羊	followed by	0.820389	灰太狼决定要抓只羊招待将到狼婆表母。灰太狼用红太狼乔装成慢羊羊造成被其捕获的假象,诱使其其它小羊与之交换,慢羊羊及时出现戳穿灰太狼的阴谋。灰太狼																		
43	灰太狼	喜羊羊	followed by	0.820389	灰太狼决定要抓只羊招待将到狼婆表母。灰太狼用红太狼乔装成慢羊羊造成被其捕获的假象,诱使其其它小羊与之交换,慢羊羊及时出现戳穿灰太狼的阴谋。灰太狼																		
45	灰太狼	喜羊羊	followed by	0.820389	灰太狼决定要抓只羊招待将到狼婆表母。灰太狼用红太狼乔装成慢羊羊造成被其捕获的假象,诱使其其它小羊与之交换,慢羊羊及时出现戳穿灰太狼的阴谋。灰太狼																		
47	沸羊羊	喜羊羊	followed by	0.818192	村长他们在兔子国到处寻找失踪的沸羊羊与喜羊羊,然而自己本身也处境危险。黑太帅已经出动兔子军队追捕小羊与沸羊羊。其实这时的沸羊羊已被囚于兔王监狱。一无																		
49	沸羊羊	喜羊羊	followed by	0.816808	灰太狼夫妇投书了黑太帅,接受了黑太帅的命令,带领了羊村的蜜蜂去围捕喜羊羊他们,就在这么危急关头,喜羊羊他们却发现了黑太帅的鬼魂,那他们真的																		
51	灰太狼	喜羊羊	followed by	0.813967	灰太狼利用装置围捕他人,羊们为了羊村的命运,纷纷出山帮助。只有懒羊羊逃过灾难,因为没看见装置,懒羊羊向灰太狼挑战,承诺只要灰太狼能捉到懒羊羊																		
53	灰太狼	喜羊羊	followed by	0.813399	灰太狼联合了武器先进的七大恶狼,偷盗了羊村并建起村庄。成功逃跑的众小羊重新武装,突破层层的重围,反攻狼堡。灰太狼以村长为头人要求快小羊,喜羊羊																		
55	慢羊羊	喜羊羊	said to be the same as	0.011025	每当慢羊羊上山采药时,总要带一剂帮手帮忙。懒羊羊用千方百计,终于得到那剂帮手。但慢羊羊却不知道,那帮手原来是灰太狼假扮的。灰太狼假扮了帮手,被灰太狼																		
57	灰太狼	喜羊羊	followed by	0.808112	众羊羊购买了喜羊羊的保险后玩了各种危险游戏,甚至在狼堡玩要娱乐。喜羊羊为保护众羊羊与灰太狼搏斗。众羊羊还不知情,在狼堡开起了舞会,被灰太狼全																		
59	喜羊羊	狼婆表母	followed by	0.808112	众羊羊购买了喜羊羊的保险后玩了各种危险游戏,甚至在狼堡玩要娱乐。喜羊羊为保护众羊羊与灰太狼搏斗。众羊羊还不知情,在狼堡开起了舞会,被灰太狼全																		

通过比较发现，从训练数据集上来看，wiki80数据集相比于TACRED数据集训练效果更好，从训练模型上看，使用BERT编码器相比于CNN编码器，在人物关系抽取上的训练效果更好。但是，通过对这5中模型结果的筛选和比较，一方面发现关系抽取的结果并没有最开始爬虫的关系抽取结果好，另一方面从关系抽取以及关系置信度上看，抽取的结果中并没有筛选到可以补充添加的人物关系。

### 3.3 基于chinese-bert-wwm的关系抽取

考虑到英文数据集对模型训练效果的影响，接下来采用20w中文人物关系数据对 chinese-bert-wwm 模型进行训练，对应项目代码放置在 `喜羊羊与灰太狼\关系抽取/OpenNRE-master` 中。

首先，下载中文 bert 模型：

```
git clone https://github.com/taorui-plus/OpenNRE.git
```

接着，在 OpenNRE-master/pretrain 目录下放置 chinese\_wwm\_pytorch 模型，其下载地址为<https://github.com/yycui/Chinese-BERT-wwm>，打开

`example/train_people_chinese_bert_softmax.py`，将代码中 'test\_chinese' 替换为 'people-relation'，并将文件打开方式中添加 `encoding='gb18030', errors='ignore'`，修改训练参数后运行，最后保存模型参数至 `ckpt/test_chinese_bert_softmax4.pth.tar`。生成中文bert预训练模型参数代码部分如下：

```
parser = argparse.ArgumentParser()
parser.add_argument('--mask_entity', action='store_true', help='Mask entity
mentions')
args = parser.parse_args()

# Some basic settings
root_path = '.'
sys.path.append(root_path)
if not os.path.exists('ckpt'):
    os.mkdir('ckpt')
ckpt = 'ckpt/test_chinese_bert_softmax.pth.tar'

# Check data
rel2id = json.load(open('E:\ZSTP\OpenNRE-
master/benchmark/test_chinese/test_chinese_rel2id.json', encoding= 'gb18030',
errors= 'ignore'))

#Define the sentence encoder
sentence_encoder = opennre.encoder.BERTEncoder(
    max_length=80,

    #pretrain_path=os.path.join(root_path, 'pretrain/chinese-bert-wwm'),
    pretrain_path='E:\ZSTP\OpenNRE-master\pretrain\chinese-bert-wwm',
    mask_entity=args.mask_entity
)

# Define the model
model = opennre.model.SoftmaxNN(sentence_encoder, len(rel2id), rel2id)

# Define the whole training framework
framework = opennre.framework.SentenceRE(
    train_path='E:\ZSTP\OpenNRE-master/benchmark/test_chinese/people-
relation_train.txt',
```

```

    val_path='E:\ZSTP\OpenNRE-master/benchmark/test_chinese/people-
relation_val.txt',
    test_path='E:\ZSTP\OpenNRE-master/benchmark/test_chinese/people-
relation_val.txt',
    model=model,
    ckpt=ckpt,
    batch_size=16, # Modify the batch size w.r.t. your device
    max_epoch=4,
    lr=2e-5,
    opt='adamw'
)

# Train the model
framework.train_model()

# Test the model
framework.load_state_dict(torch.load(ckpt)['state_dict'])
result = framework.eval_model(framework.test_loader)

# Print the result
print('Accuracy on test set: {}'.format(result['acc']))

```

训练结果如下:

```

=== Epoch 0 train ===
100%|██████████| 12375/12375 [1:29:22<00:00, 2.31it/s, acc=0.559, loss=1.33]
=== Epoch 0 val ===
100%|██████████| 63/63 [00:10<00:00, 5.88it/s, acc=0.704]
Best ckpt and saved.
=== Epoch 1 train ===
100%|██████████| 12375/12375 [1:20:33<00:00, 2.56it/s, acc=0.704, loss=0.87]
=== Epoch 1 val ===
100%|██████████| 63/63 [00:10<00:00, 5.89it/s, acc=0.823]
Best ckpt and saved.
=== Epoch 2 train ===
100%|██████████| 12375/12375 [1:20:30<00:00, 2.56it/s, acc=0.791, loss=0.607]
=== Epoch 2 val ===
100%|██████████| 63/63 [00:10<00:00, 5.86it/s, acc=0.886]
Best ckpt and saved.
=== Epoch 3 train ===
100%|██████████| 12375/12375 [1:20:29<00:00, 2.56it/s, acc=0.847, loss=0.442]
=== Epoch 3 val ===
100%|██████████| 63/63 [00:10<00:00, 5.87it/s, acc=0.919]
Best ckpt and saved.
Best acc on val set: 0.919000
100%|██████████| 63/63 [00:10<00:00, 5.88it/s, acc=0.919]
Accuracy on test set: 0.919

Process finished with exit code 0

```

最终 chinese-bert-wm 模型ACC达到0.919。

将模型预训练参数 test\_chinese\_bert\_softmax.pth.tar 加载至 chinese-bert-wwm 的 framework.model 中, 对实验数据进行关系抽取, 并将预测结果按照 ["人物A", "人物B", "关系", "概率", "文本"] 格式输出至 ans\_chinese\_bert\_wwm.csv 中。对应代码在 喜羊羊与灰太狼\关系抽取\OpenNRE-master\example\predict\_xyy\_chinese\_bert.py 文件中。

```
# encoding:gbk

parser = argparse.ArgumentParser()
parser.add_argument('--mask_entity', action='store_true', help='Mask entity
mentions')
args = parser.parse_args()

# Some basic settings
root_path = '.'
sys.path.append(root_path)
if not os.path.exists('ckpt'):
    os.mkdir('ckpt')
ckpt = 'ckpt/test_chinese_bert_softmax3.pth.tar'

# Check data
rel2id = json.load(open('OpenNRE-
master/benchmark/test_chinese/test_chinese_rel2id.json', encoding= 'gb18030',
errors= 'ignore'))

#Define the sentence encoder
sentence_encoder = opennre.encoder.BERTEncoder(
    max_length=80,

    #pretrain_path=os.path.join(root_path, 'pretrain/chinese-bert-wwm'),
    pretrain_path='OpenNRE-master\pretrain\chinese-bert-wwm',
    mask_entity=args.mask_entity
)

# Define the model
model = opennre.model.SoftmaxNN(sentence_encoder, len(rel2id), rel2id).cuda()

framework = opennre.framework.SentenceRE(
    train_path='OpenNRE-master/benchmark/test_chinese/test_chinese_train.txt',
    val_path='OpenNRE-master/benchmark/test_chinese/test_chinese_val.txt',
    test_path='OpenNRE-master/benchmark/test_chinese/test_chinese_val.txt',
    model=model,
    ckpt=ckpt,
    batch_size=16, # Modify the batch size w.r.t. your device
    max_epoch=7,
    lr=2e-5,
    opt='adamw'
)
framework.load_state_dict(torch.load(ckpt)['state_dict'])
import csv
fw=open('ans_chinese_bert_wwm.csv', "w")
csv_writer = csv.writer(fw)
csv_writer.writerow(["人物A", "人物B", "关系", "概率", "文本"])
with open('textth.txt', 'r', encoding='gb18030', errors='ignore') as f :
    for line in f :
        if line[0].isdigit() :
            continue
```



```

elif len(line) != 0 :
    ss=eval(line)
    #print(ss['text'][ss['h']['pos'][0]:ss['h']['pos'][1]+1], ss['text']
    [ss['t']['pos'][0]:ss['t']['pos'][1]+1])
    result = framework.model.infer(ss)
    relation=result[0].encode('gb18030').decode('utf-8')
    #print(relation)
    a=ss['text'][ss['h']['pos'][0]:ss['h']['pos'][1]+1]
    b=ss['text'][ss['t']['pos'][0]:ss['t']['pos'][1]+1]
    if a==b:
        continue
    else:
        csv_writer.writerow([a, b, relation,result[1],ss['text']])
fw.close()

```

ans\_chinese\_bert\_wwm.csv 输出结果如下:

从结果上看, 关系抽取结果存在以下两个问题:

- 第一, 由于文本本身数据清洗不全面, 导致“灰太”, “招潮蟹”等词语滞留在输出文件中。
- 第二, 由于模型训练数据的关系是:

```

{"父母":0, "夫妻":1, "师生":2, "兄弟姐妹":3, "合作":4, "情侣":5, "祖孙":6, "好友":7, "亲戚":8, "同门":9, "上下级":10, "unknown":11}

```

即模型训练时未涉及“敌人”关系, 导致将一些敌对关系, 例如“黑大帅”与“兔王”, “老虎”与“灰太狼”的关系识别为“好友”关系。在这里, 为了更全面展现人物之间的关系连接, 本实验场景将“敌人”关系划为“好友”关系。接着, 将有效关系的关系置信区间设置为0.8, 发现符合的输出结果有38个, 然后对原本爬虫得到的 data.csv 中没有的三元组关系进行人为补充, 并将“好友”关系的描述统一为“朋友”。共添加了4条关系三元组, 补充结果如下:

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V
1451	黑牛国王	黑细菌	物种																			
1452	黑牛国王	Blacky	英文名																			
1453	黑牛国王	毛方圆	(配音演员)																			
1454	黑钻	科技城	出身地区																			
1455	黑钻	黑钻	别名																			
1456	黑钻	小雨	妹妹																			
1457	黑钻	男	性别																			
1458	黑钻	小雷	本名																			
1459	黑钻	科技城	活动范围																			
1460	黑钻	豹	物种																			
1461	黑钻	赵娜	配音演员																			
1462	鼠一	鼠一	本名																			
1463	鼠一	鼠	物种																			
1464	鼠一	刘红韵	配音演员																			
1465	黑大帅	兔王	朋友																			
1466	黑大帅	黑大狼	朋友																			
1467	灰太狼	黑大帅	朋友																			
1468	老虎	灰太狼	朋友																			
1469																						
1470																						
1471																						
1472																						
1473																						
1474																						
1475																						
1476																						
1477																						
1478																						

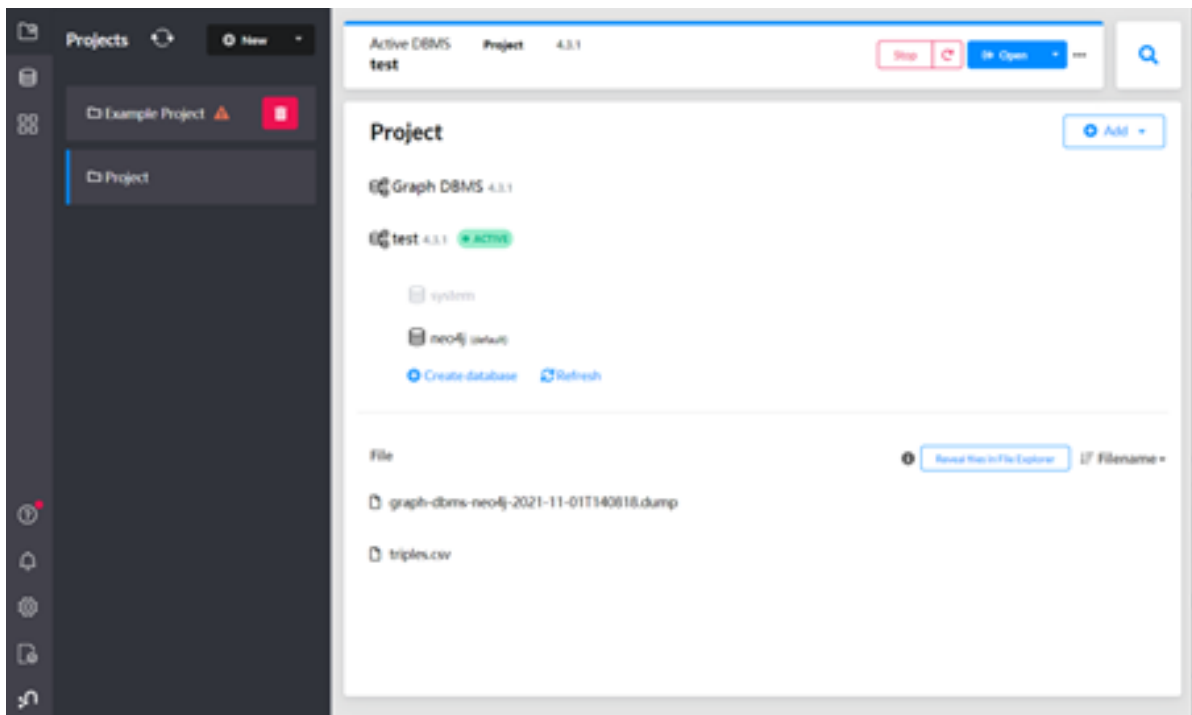
最后，基于6种模型的关系抽取训练与推理完毕。从实验结果上看，基于 chinese-bert-www 模型的关系抽取相比于另外5种模型，准确率有较为明显的提升，同时抽取结果更符合实验文本的语义场景。

## 4 可视化与知识问答

通过 Neo4j Deskto，对《喜羊羊与灰太狼》的关系抽取结果进行可视化呈现，并实现QA。对应项目代码在 [喜羊羊与灰太狼\可视化与知识问答](#) 目录下。

### 4.1 可视化

首先，下载 Neo4j Deskto，解压安装，创建 test 数据库并激活运行：



接着，将 relation.txt 三元组导入 Neo4j Browser，进行可视化，具体代码在 [喜羊羊与灰太狼\可视化与知识问答/neo.py](#) 目录下。

```
class NeoDB(object):
    def __init__(self):
```

```

self.graph = Graph("http://localhost:7474",
                   user="neo4j", password="test")
self.CA_LIST = {"青青草原": 0, "青青草原": 1,
                "远古时代": 2, "龙世界": 3, "海底": 4, "其他": 5, "冰极": 6}
self.similar_words = {
    "爸爸": "父亲", "妈妈": "母亲", "爸": "父亲", "妈": "母亲", "朋友": "好朋
友",
    "哥": "哥哥"
}

def create_graph(self, relation_filename="data/relation.txt"):
    g = Graph('http://localhost:7474',user='neo4j',password='test')
    with open(relation_filename,'r',encoding='utf-8') as f:
        for line in f.readlines():
            rela_array = line.strip("\n").split(",")
            print(rela_array)
            self.graph.run("MERGE(p: Person{cate:'%s',Name: '%s'})" % (
                rela_array[3], rela_array[0]))
            self.graph.run("MERGE(p: Person{cate:'%s',Name: '%s'})" % (
                rela_array[4], rela_array[1]))
            self.graph.run(
                "MATCH(e: Person), (cc: Person) \
                WHERE e.Name='%s' AND cc.Name='%s'\
                CREATE(e)-[r:%s{relation: '%s'}]->(cc)\
                RETURN r" % (rela_array[0], rela_array[1], rela_array[2],
rela_array[2]))
        )
    return

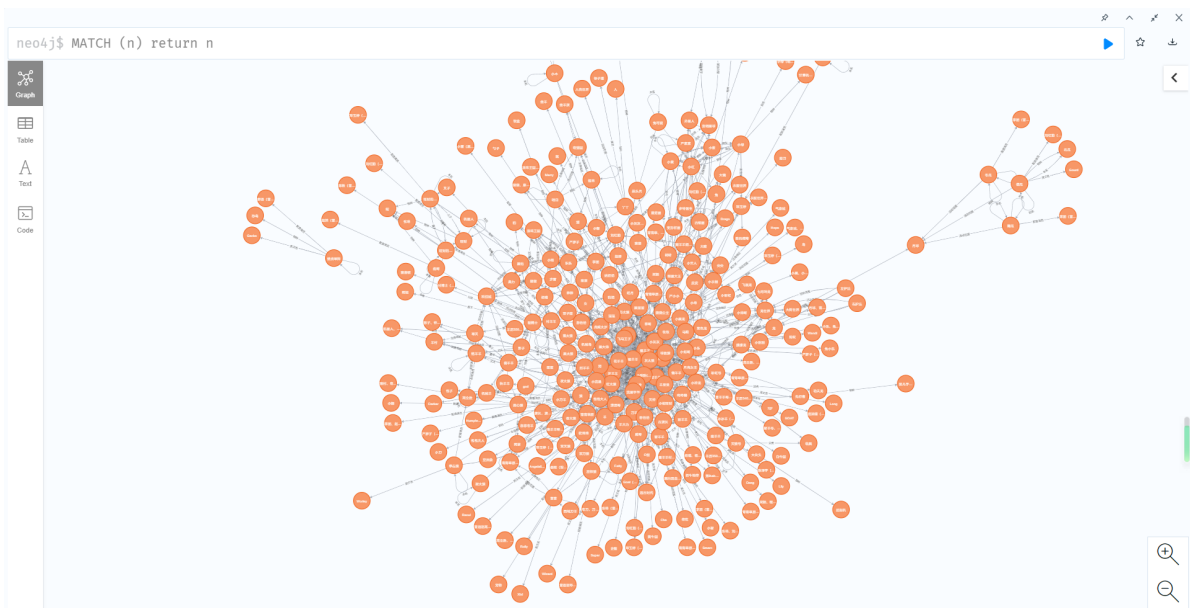
if __name__ == "__main__":
    neo_db = NeoDB()
    neo_db.create_graph("./data/relation.txt")

```

最后，(1) 可视化整张图谱。在 Neo4j Browser 端输入：

```
MATCH (n) retu
```

完整图谱展现结果如下：









```

        return '问题有误'
    try:
        data = self.neo.graph.run(query)
        data = list(data)[0]
        logging.debug(str(data))

        result = data['n.cate']+'的'+data['n.Name']
        for item in words:
            result += '的'
            result += item
            result += '是'
            result += data['p.cate']+'的'+data['p.Name']
        return result
    except Exception as e:
        return '没有答案'

```

例如，输入问题：

**用户：小灰灰的父亲是谁？**

解析得到，`小灰灰` 和 `父亲` 两个词语。

接着，采用Neo4j查询语句查询知识图谱，返回问题答案。

例如，输入问题：`小灰灰的爸爸是？`，返回回答为：`灰太狼`。

**用户：小灰灰的爸爸是？**

**回答：灰太狼**

实践其他问答结果如下：

**用户：小灰灰的爸爸是？**

**回答：青青草原的小灰灰的爸爸是狼堡的灰太狼**

**用户：小灰灰的妈妈是？**

**回答：青青草原的小灰灰的爸爸是狼堡的红太狼**

**用户：红太狼的妈妈是？**

**回答：狼堡的的红太狼的妈妈是狼堡的粉红太狼**

**用户：喜羊羊的爸爸是？**

**回答：青青草原的的喜羊羊的爸爸是青青草原的智羊羊**

**用户：喜羊羊的妈妈是？**

**回答：青青草原的的喜羊羊的爸爸是青青草原的丽羊羊**

**用户：智羊羊的妻子是？**

**回答：青青草原的的智羊羊的妻子是青青草原的丽羊羊**

**用户：灰太狼的爸爸是？**

**回答：狼堡的的灰太狼的爸爸是狼堡的黑太狼**

