

第三期『千言万语』系列技术分享

篇章事件抽取中的要素组合问题

分享嘉宾：朱桐 苏州大学博士生

信息抽取 (Information Extraction) 中的几个关键任务



.....

事件抽取 Event Extraction

事件检测
Event Detection

要素抽取
Argument Extraction

触发词识别
Trigger Identification

触发词分类
Trigger Classification

要素识别
Argument Identification

要素分类
Argument Classification

命名实体识别 Named Entity Recognition

实体分类
Entity Typing

实体识别
Entity Identification

实体关系抽取 Entity Relationship Extraction

关系识别
Relation Identification

关系分类
Relation Classification



什么是事件抽取?

- 将非结构化文本转为结构化数据的过程

海绵宝宝今天晚上8点要在章鱼哥家里给蟹老板过生日。





事件模板与要素组合

| 生日宴会 | |
|------|--------|
| 时间 | 今天晚上8点 |
| 地点 | 章鱼哥的家 |
| 寿星 | 蟹老板 |
| 参与人 | 海绵宝宝 |
| 主菜 | N/A |

事件实例
Event Record

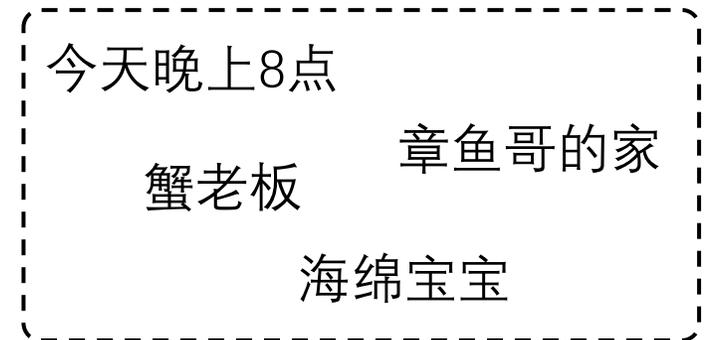
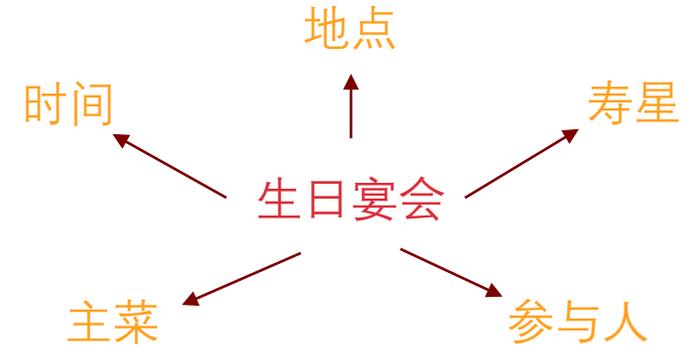


事件模板
Event Template

所有类型事件模板
的集合被称为 schema

事件要素组合
Event Argument Combination

一个事件实例中所有要素
所构成的集合





事件要素组合问题

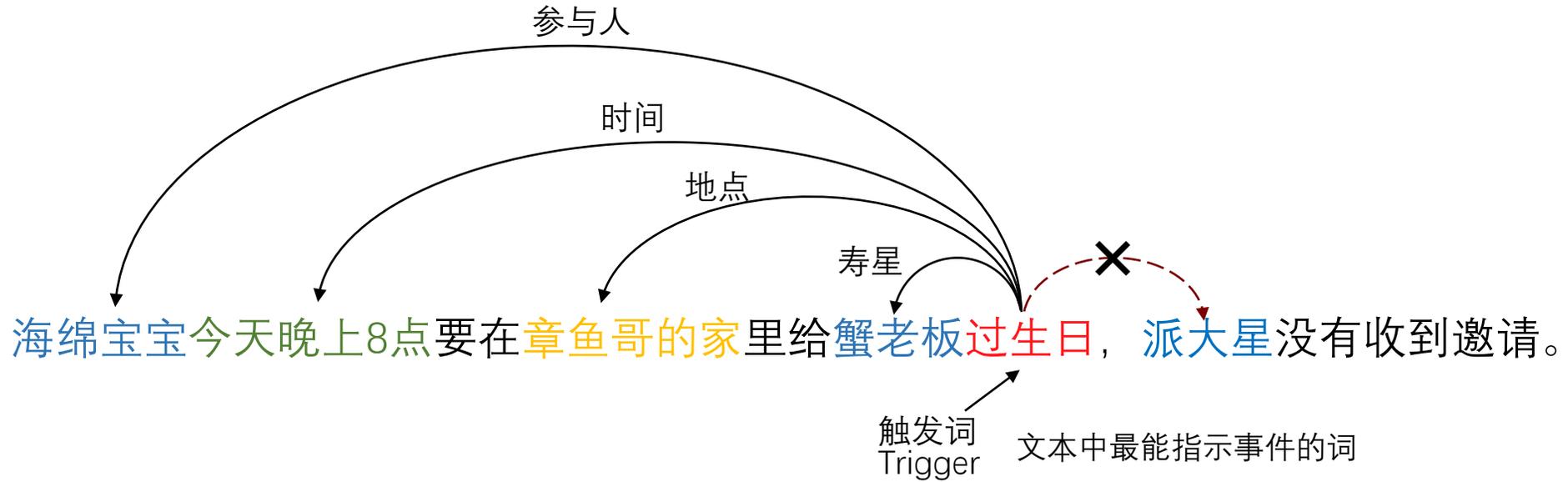
海绵宝宝今天晚上8点要在章鱼哥家里给蟹老板过生日，派大星没有收到邀请。

共有 2^5 种要素组合情况





传统带触发词的要素组合方案





篇章事件抽取的特点

海绵宝宝早就等着这一天了。

今年1月份的时候，珍妮告诉他，自己的爸爸要过生日了，希望能给他办一个惊喜派对。

海绵宝宝答应了她，并和好朋友珊迪说了这件事情。

珊迪说他们应该叫更多的朋友一起，比如章鱼哥和派大星。

大家商量了一下，决定今天晚上8点在章鱼哥的家里举行生日派对。

而此时的蟹老板还不知道等待他的是怎样的一个惊喜。

- 1、文本很长
 - 2、实体在文中出现多次，有多个提及 (mention)
 - 3、事件参与者通常分散在句子的各个位置
- } 平均信息密度比较低，更符合实际应用



标注困难，触发词缺失



触发词标注问题1：无标注

"证券代码：300142证券简称：沃森生物公告编号：2016-072"，
"云南沃森生物技术股份有限公司关于股东解除股权质押的公告"，
"本公司及董事会全体成员保证信息披露内容的真实、准确和完整，没有虚假记载、误导性陈述或重大遗漏。"，
"云南沃森生物技术股份有限公司（以下简称“公司”）日前接到股东李云春先生函告，获悉李云春先生所持有的本公司部分股份解除质押，具体情况如下："，
"李云春先生曾于2015年5月7日同招商证券股份有限公司就其持有的13596398股公司股票办理了股票质押回购业务，质押期限自2015年5月7日起"，
"至质权人向中国证券登记结算有限责任公司深圳分公司办理解除质押登记为止（详见公司在证监会指定的信息披露网站巨潮资讯网披露的第2015-048号公告）。"，
"李云春先生于2016年5月6日在中国证券登记结算有限责任公司深圳分公司办理了40789194股（含除权后派送的27192796股股份）公司股份的质押解除手续。"，
"李云春先生本次解除质押的公司股份占公司股份总数的2.91%，占其所持公司股份的25.16%。"，
"截至本公告披露日，李云春先生共持有公司股份162103218股，占公司股份总数的11.55%。"，
"李云春先生共质押其持有的公司股份86304393股，占公司股份总数的6.15%，占其所持公司股份的53.24%。"，
"特此公告。"，
"云南沃森生物技术股份有限公司"，
"董事会"，
"二〇一六年五月九日"

ChFinAnn数据集基于远程监督构建，其中完全不提供任何触发词，只有事件要素标注

```
"EquityPledge",  
{  
  "Pledger": "李云春",  
  "PledgedShares": "40789194股",  
  "Pledgee": "招商证券股份有限公司",  
  "TotalHoldingShares": "162103218股",  
  "TotalHoldingRatio": "11.55%",  
  "TotalPledgedShares": "86304393股",  
  "StartDate": "2015年5月7日",  
  "EndDate": "2016年5月6日",  
  "ReleasedDate": "2016年5月6日"
```



触发词标注问题1：触发词被不同事件实例共用

```
{
  "id": "e55ffdd8af2e77df58b3020c8ee5502a",
  "title": "苹果已收购企业设备管理公司Fleetsmith",
  "text": "原标题: 苹果已收购企业设备管理公司Fleetsmith 来源: 威锋网\n苹果已经收购了企业设备管理公司 Fleetsmith, 该公司于周三宣布了这一消息.\n成立于 2014 年的 Fleetsmith 在博客文章中宣布他们现在已成为苹果的一员, 并称他们“很高兴加入苹果”。该公司为 IT 部门提供企业解决方案, 以管理 Mac, iPad 和 iPhone.\n该公司写道: “我们的共同价值观是在不牺牲隐私和安全的条件下将客户放在我们所做一切工作的中心, 这意味着我们可以真正实现我们的使命, 将 Fleetsmith 提供给世界各地各种规模的企业和机构。”\n在收购的同一周, 苹果在 WWDC 2020 举行了“管理苹果设备的新功能”开发人员会议。在会议期间, 苹果宣布了 Mac Pro 的新管理功能, Mac Supervision 的更改以及 macOS Big Sur 中的托管软件更新及其它功能.\n通过 Fleetsmith 的收购, 苹果可能希望进一步为企业和教育客户增强其第一方设备管理选项。到目前为止, 苹果公司主要依靠第三方解决方案为其客户提供 MDM 平台.\nFleetsmith 是苹果在 2020 年的又一笔收购, 之前的收购包括天气应用 Dark Sky, 实时 VR 活动公司 NextVR 和 AI 初创公司 Voysis.",
  "event_list": [
    {
      "trigger": "收购",
      "event_type": "企业收购",
      "arguments": [
        {"role": "收购方", "argument": "苹果"},
        {"role": "被收购方", "argument": "Fleetsmith"},
        {"role": "披露时间", "argument": "周三"}
      ]
    },
    {
      "trigger": "收购",
      "event_type": "企业收购",
      "arguments": [
        {"role": "收购方", "argument": "苹果"},
        {"role": "被收购方", "argument": "Dark Sky"},
        {"role": "被收购方", "argument": "NextVR"},
        {"role": "被收购方", "argument": "Voysis"},
        {"role": "收购完成时间", "argument": "2020 年"}
      ]
    }
  ]
}
```

DuEE-Fin数据集中提供触发词，但无位置标注，因此触发词可能被不同实例共用



篇章事件抽取的特点

海绵宝宝早就等着这一天了。
今年1月份的时候，珍妮告诉他，自己的爸爸要过生日了，希望能给他办一个惊喜派对。
海绵宝宝答应了她，并和好朋友珊迪说了这件事情。
珊迪说他们应该叫更多的朋友一起，比如章鱼哥和派大星。
大家商量了一下，决定今天晚上8点在章鱼哥的家里举行生日派对。
而此时的蟹老板还不知道等待他的是怎样的一个惊喜。



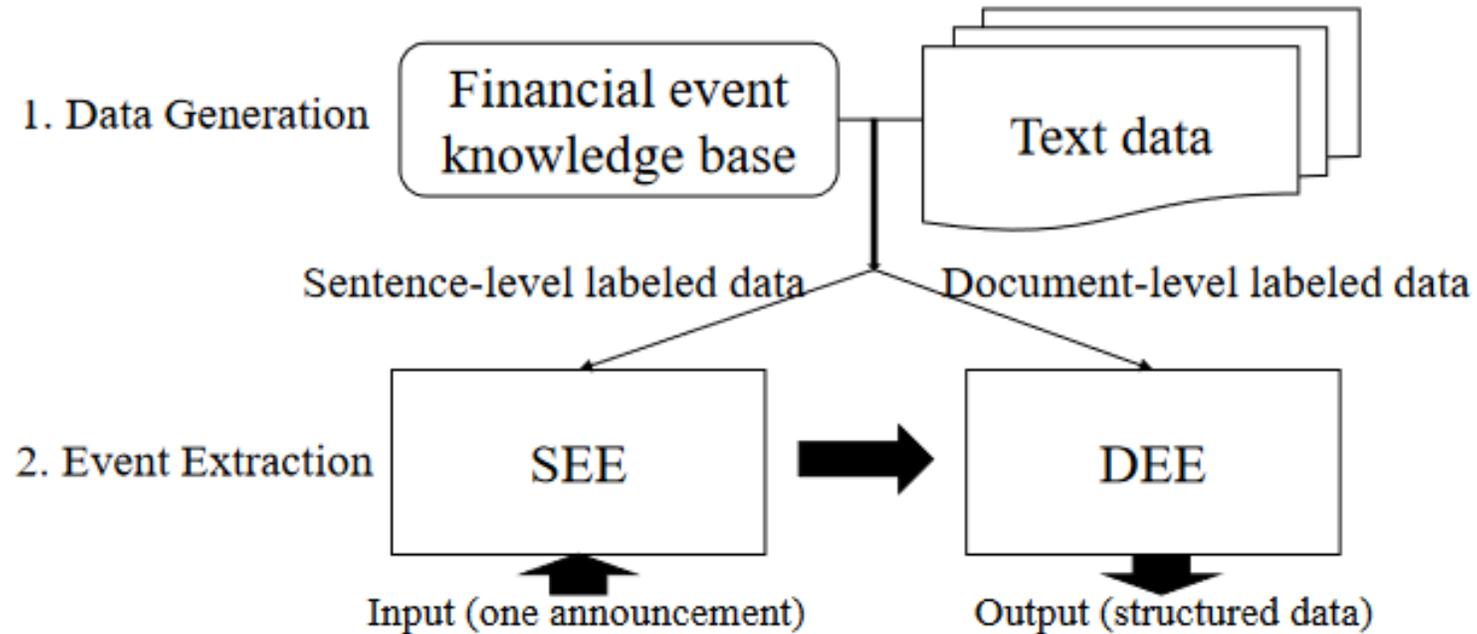
标注困难，触发词缺失、共用、质量差等情况很常见



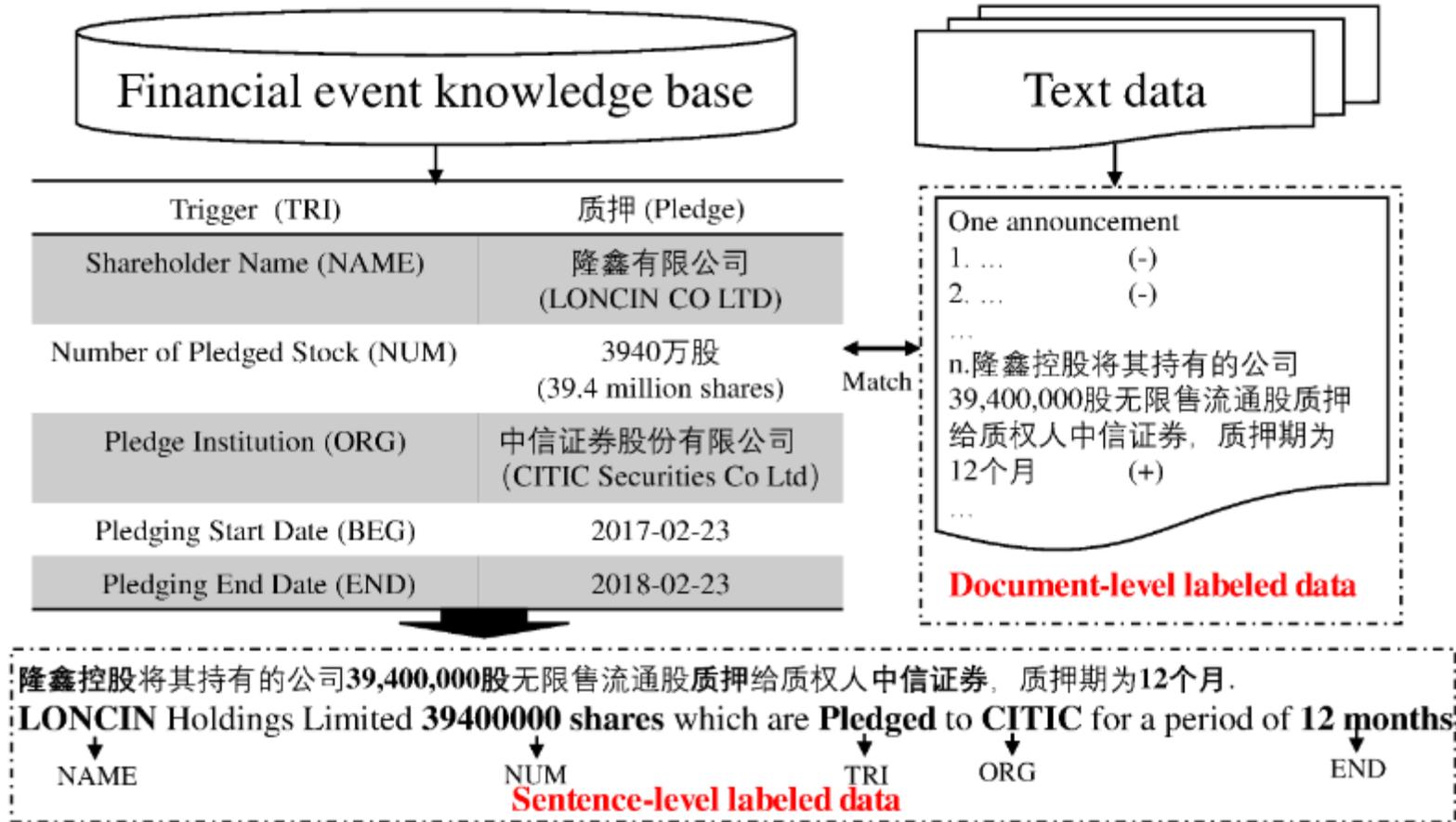
怎么进行要素的组合，以完成事件抽取呢？

DCFEE: 整体方案

- 构建了中文篇章事件抽取的数据集
- 提供了一种篇章事件抽取方法

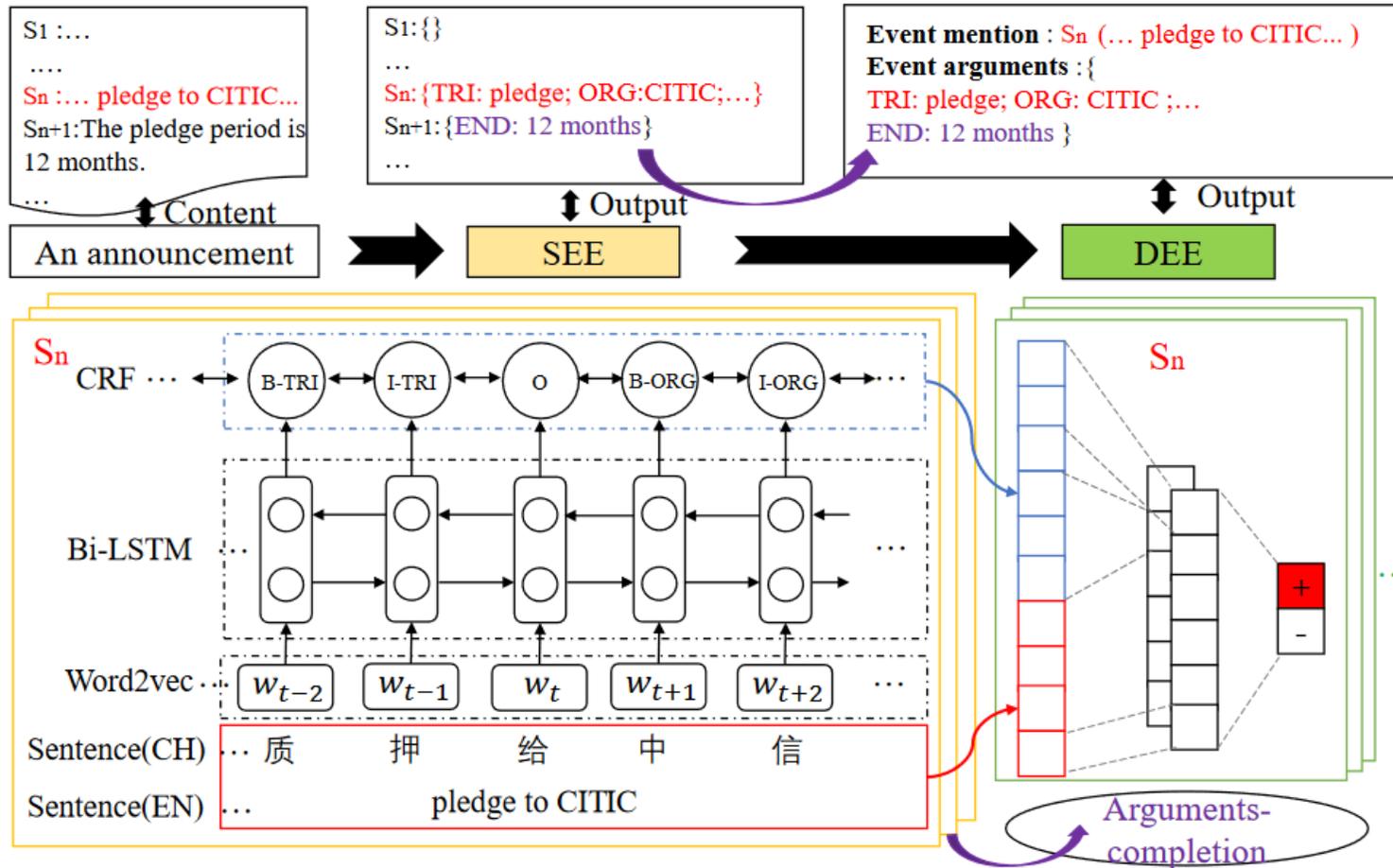


DCFEE: 数据构建



- 使用金融领域知识图谱和金融公告做**远程监督** (Distant Supervision, DS)
- 选择要素匹配最多的句子作为**关键句**, 并从触发词词典中匹配该关键句中的触发词

DCFEE: 具体方法



- 首先使用LSTM+CRF的方法抽取篇章中的所有要素和触发词
- 关键句检测: 之后对每个句子做分类, 判断当前句子是否为**关键句**
- 要素补全: 如果事件实例不完整, 则在关键句周围的其它句子中寻找其它要素



DCFEE: 例子

1. 抽取所有的要素和触发词

海绵宝宝早就等着这一天了。

今年1月份的时候，珍妮告诉他，自己的爸爸要过生日了，希望能给他办一个惊喜派对。

海绵宝宝答应了她，并和好朋友珊迪说了这件事情。

珊迪说他们应该叫更多的朋友一起，比如章鱼哥和派大星。

大家商量了一下，决定今天晚上8点在章鱼哥的家里举行生日派对。

而此时的蟹老板还不知道等待他的是怎样的一个惊喜。



DCFEE: 例子

1. 抽取所有的要素和触发词

海绵宝宝早就等着这一天了。

今年1月份的时候，珍妮告诉他，自己的爸爸要过生日了，希望能给他办一个惊喜派对。

海绵宝宝答应了她，并和好朋友珊迪说了这件事情。

珊迪说他们应该叫更多的朋友一起，比如章鱼哥和派大星。

大家商量了一下，决定今天晚上8点在章鱼哥的家里举行生日派对。←—— 2. 找到关键句

而此时的蟹老板还不知道等待他的是怎样的一个惊喜。

| 生日宴会 | |
|------|--------|
| 时间 | 今天晚上8点 |
| 地点 | 章鱼哥的家 |



DCFEE: 例子

1. 抽取所有的要素和触发词

海绵宝宝早就等着这一天了。

今年1月份的时候，珍妮告诉他，自己的爸爸要过生日了，希望能给他办一个惊喜派对。

海绵宝宝答应了她，并和好朋友珊迪说了这件事情。

珊迪说他们应该叫更多的朋友一起，比如章鱼哥和派大星。

大家商量了一下，决定今天晚上8点在章鱼哥的家里举行生日派对。

而此时的蟹老板还不知道等待他的是怎样的一个惊喜。

2. 找到关键句

3. 在关键句周围寻找更多的要素补全事件

| 生日宴会 | |
|------|------------|
| 时间 | 今天晚上8点 |
| 地点 | 章鱼哥的家 |
| 参与人 | 珊迪、章鱼哥、派大星 |
| 寿星 | 蟹老板 |



DCFEE: 例子

1. 抽取所有的要素和触发词

海绵宝宝早就等着这一天了。

今年1月份的时候，珍妮告诉他，自己的爸爸要过生日了，希望能给他办一个惊喜派对。

海绵宝宝答应了她，并和好朋友珊迪说了这件事情。

珊迪说他们应该叫更多的朋友一起，比如章鱼哥和派大星。

大家商量了一下，决定今天晚上8点在章鱼哥的家里举行生日派对。

而此时的蟹老板还不知道等待他的是怎样的一个惊喜。

2. 找到关键句

3. 在关键句周围寻找更多的要素补全事件

| 生日宴会 | |
|------|------------|
| 时间 | 今天晚上8点 |
| 地点 | 章鱼哥的家 |
| 参与人 | 珊迪、章鱼哥、派大星 |
| 寿星 | 蟹老板 |

- 1、窗口的选择对结果影响很大
- 2、全局信息利用不足，容易遗漏要素
- 3、在含有多个事件的复杂文本中，可能会产生很多FP错误



| Stage | SEE | | | DEE | | |
|-----------|---------|---------|-----------|---------|---------|-----------|
| Type | $P(\%)$ | $R(\%)$ | $F_1(\%)$ | $P(\%)$ | $R(\%)$ | $F_1(\%)$ |
| <i>EF</i> | 90.00 | 90.41 | 90.21 | 80.70 | 63.40 | 71.01 |
| <i>EP</i> | 93.31 | 94.36 | 93.84 | 80.36 | 65.91 | 72.30 |
| <i>ER</i> | 92.79 | 93.80 | 93.29 | 88.79 | 82.02 | 85.26 |
| <i>EO</i> | 88.76 | 91.88 | 90.25 | 80.77 | 45.93 | 58.56 |

Table 4: P , R , F_1 of SEE, DEE on the different event types.

- 句子级事件抽取的效果很好，篇章事件抽取效果欠佳，特别是召回率很低



多实例的情况

增加了新的
“生日聚会”
事件和“出
售”事件

海绵宝宝早就等着这一天了。

今年1月份的时候，珍妮告诉他，自己的爸爸要过生日了，希望能给他办一个惊喜派对。

去年的今天，泡芙小姐给蟹老板办了一场惊喜派对，他非常开心。

但是珍妮手头没有足够的钱，希望海绵宝宝能够帮帮她。

海绵宝宝答应了她，并和好朋友珊迪说了这件事情。

珊迪和海绵宝宝在星期天找到了痞老板，并把蟹黄堡配方以500元的价格卖给了他。

这下他们有足够的钱啦~

珊迪说他们应该叫更多的朋友一起参加，比如章鱼哥和派大星。

大家商量了一下，决定今天晚上8点在章鱼哥家里举行生日派对。

而此时的蟹老板还不知道等待他的是怎样的一个惊喜。

- 在无触发词的多实例场景下，我们应该怎么办？

| Entity Mark Table | | |
|-------------------|------------|-----------------------------|
| Mark | Entity | Entity (English) |
| [PER] | 刘维群 | Weiqun Liu |
| [ORG] | 国信证券股份有限公司 | Guosen Securities Co., Ltd. |
| [DATE1] | 2017年9月22日 | Sept. 22nd, 2017 |
| [DATE2] | 2018年9月6日 | Sept. 6th, 2018 |
| [DATE3] | 2018年9月20日 | Sept. 20th, 2018 |
| [DATE4] | 2019年3月20日 | Mar. 20th, 2019 |
| [SHARE1] | 750000股 | 750000 shares |
| [SHARE2] | 975000股 | 975000 shares |
| [SHARE3] | 525000股 | 525000 shares |
| [SHARE4] | 1500000股 | 1500000 shares |
| [SHARE5] | 16768903股 | 16768903 shares |
| [RATIO] | 1.0858% | 1.0858% |

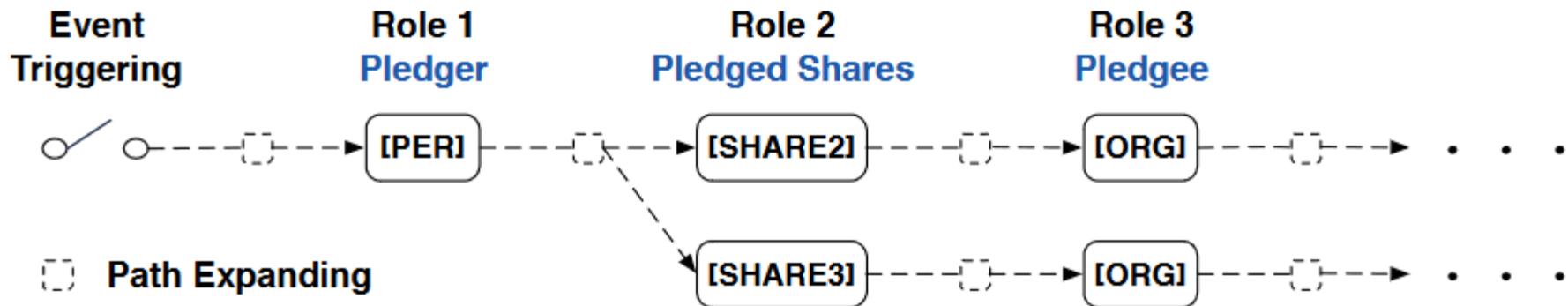
| Event Table of Equity Pledge | | | | | | |
|------------------------------|----------------|---------|------------|----------|----------------------|---------------------|
| Pledger | Pledged Shares | Pledgee | Begin Date | End Date | Total Holding Shares | Total Holding Ratio |
| [PER] | [SHARE2] | [ORG] | [DATE1] | [DATE4] | [SHARE5] | [RATIO] |
| [PER] | [SHARE3] | [ORG] | [DATE2] | [DATE4] | [SHARE5] | [RATIO] |

| Event Role | Event Record | Event Argument | Entity Mention |
|------------|--------------|----------------|--|
| | | ID | Sentence |
| | | 5 | [DATE1], [PER]将其持有的公司[SHARE1]股份质押给[ORG]. In [DATE1], [PER] pledged his [SHARE1] to [ORG]. |
| | | 7 | 公司实施资本公积金转增股本后，其质押股份变为[SHARE2]. After the company carried out the transferring of the capital accumulation fund to the capital stock, his pledged shares became [SHARE2]. |
| | | 8 | [DATE2], [PER]将其持有的[SHARE3]公司股份质押给[ORG]，作为对上述质押股份的补充质押。 In [DATE2], [PER] pledged [SHARE3] to [ORG], as a supplementary pledge to the above pledged shares. |
| | | 9 | 上述质押及补充质押股份合计为[SHARE4]，原定购回日期为[DATE3]. The aforementioned pledged and supplementary pledged shares added up to [SHARE4], and the original repurchase date was [DATE3]. |
| | | 10 | [DATE3], [PER]针对其质押的[SHARE4]股份办理了延期购回业务，购回日期延长至[DATE4]. In [DATE3], [PER] extended the repurchase date to [DATE4] for [SHARE4] he pledged. |
| | | 12 | 截至本公告日，[PER]持有公司股份[SHARE5]，占公司总股本的[RATIO]. As of the date of this announcement, [PER] hold [SHARE5] of the company, accounting for [RATIO] of the total share capital of the company. |



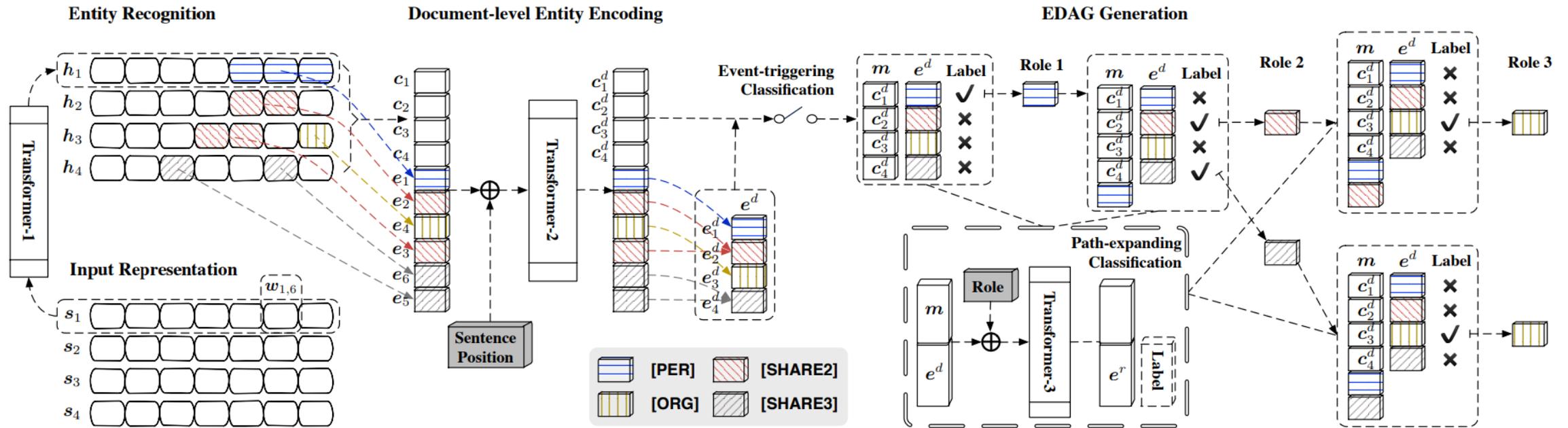
Doc2EDAG: 有向无环图进行要素组合的方案

基于实体的有向无环图 (Entity-based Directed Acyclic Graph, EDAG)



- 首先对篇章进行分类，得到事件类别
- 针对每个确定的类别，依次判断某一角色对应哪些实体
- 对于一个角色对应k个实体的情况，则将路径**分割成k条**，并继续判断下一角色对应的要素

Doc2EDAG: 整体架构



- 在训练的过程中**动态生成EDAG**
- 在每步判断要素角色对应的要素时，都要**对所有的实体进行分类**
- 所有句子的表示和当前路径的所有前序节点作为**记忆单元**增强实体表示



1. 抽取所有的要素

海绵宝宝早就等着这一天了。

今年1月份的时候，珍妮告诉他，自己的爸爸要过生日了，希望能给他办一个惊喜派对。

去年的今天，泡芙小姐给蟹老板办了一场惊喜派对，他非常开心。

但是珍妮手头没有足够的钱，希望海绵宝宝能够帮帮她。

海绵宝宝答应了她，并和好朋友珊迪说了这件事情。

珊迪和海绵宝宝在星期天找到了痞老板，并把蟹黄堡配方以500元的价格卖给了他。

这下他们有足够的钱啦~

珊迪说他们应该叫更多的朋友一起参加，比如章鱼哥和派大星。

大家商量了一下，决定今天晚上8点在章鱼哥的家里举行生日派对。

而此时的蟹老板还不知道等待他的是怎样的一个惊喜。



1. 抽取所有的要素

海绵宝宝早就等着这一天了。
今年1月份的时候, 珍妮告诉他, 自己的爸爸要过生日了, 希望能给他办一个惊喜派对。
去年的今天, 泡芙小姐给蟹老板办了一场惊喜派对, 他非常开心。
但是珍妮手头没有足够的钱, 希望海绵宝宝能够帮帮她。
海绵宝宝答应了她, 并和好朋友珊迪说了这件事情。
珊迪和海绵宝宝在星期天找到了痞老板, 并把蟹黄堡配方以500元的价格卖给了他。
这下他们有足够的钱啦~
珊迪说他们应该叫更多的朋友一起参加, 比如章鱼哥和派大星。
大家商量了一下, 决定今天晚上8点在章鱼哥的家里举行生日派对。
而此时的蟹老板还不知道等待他的是怎样的一个惊喜。

2. 事件分类

生日聚会、出售

Doc2EDAG: 例子

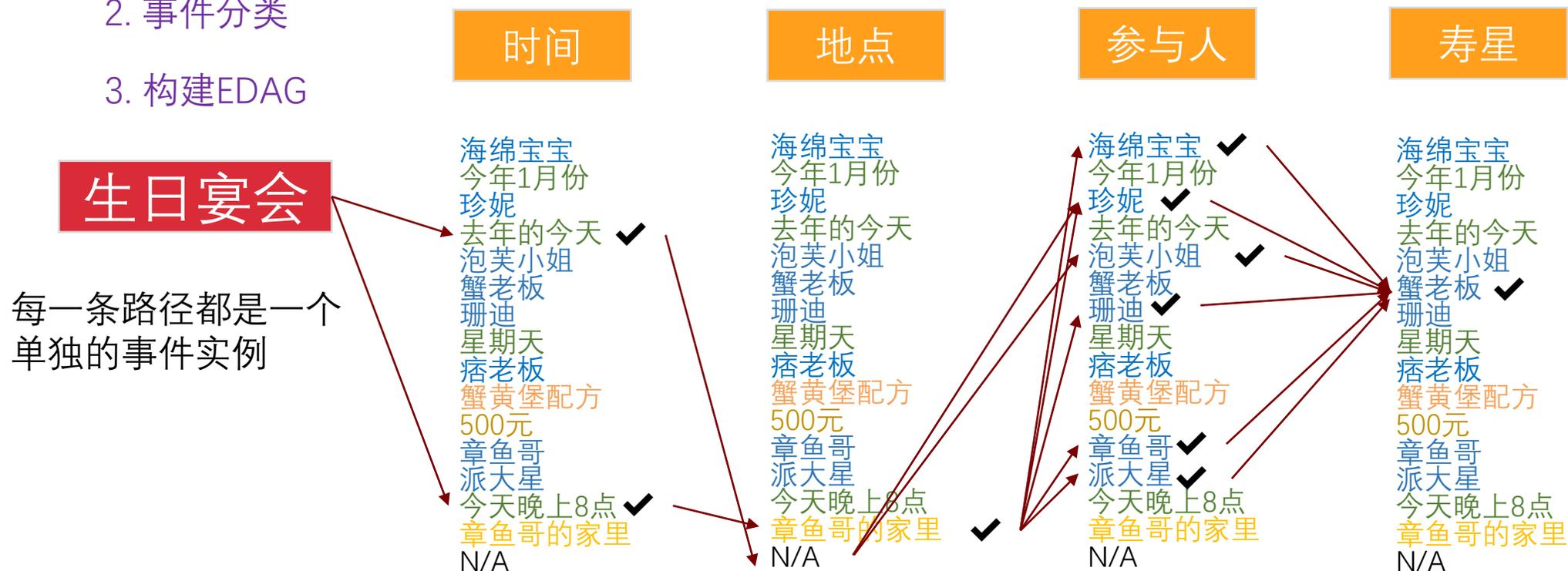


1. 抽取所有的要素

海绵宝宝早就等着这一天了。
今年1月份的时候，珍妮告诉他，自己的爸爸要过生日了，希望能给他办一个惊喜派对。
去年的今天，泡芙小姐给蟹老板办了一场惊喜派对，他非常开心。
但是珍妮手头没有足够的钱，希望海绵宝宝能够帮帮她。
海绵宝宝答应了她，并和好朋友珊迪说了这件事情。
珊迪和海绵宝宝在星期天找到了痞老板，并把蟹黄堡配方以500元的价格卖给了他。
这下他们有足够的钱啦~
珊迪说他们应该叫更多的朋友一起参加，比如章鱼哥和派大星。
大家商量了一下，决定今天晚上8点在章鱼哥的家举行生日派对。
而此时的蟹老板还不知道等待他的是怎样的一个惊喜。

2. 事件分类

3. 构建EDAG





Doc2EDAG: ChFinAnn数据集

- ChFinAnn: 使用远程监督方法自动生成的大规模数据集，且质量不错

| Event | #Train | #Dev | #Test | #Total | MER (%) |
|-------|--------|-------|-------|--------|---------|
| EF | 806 | 186 | 204 | 1,196 | 32.0 |
| ER | 1,862 | 297 | 282 | 3,677 | 16.1 |
| EU | 5,268 | 677 | 346 | 5,847 | 24.3 |
| EO | 5,101 | 570 | 1,138 | 6,017 | 28.0 |
| EP | 12,857 | 1,491 | 1,254 | 15,602 | 35.4 |
| All | 25,632 | 3,204 | 3,204 | 32,040 | 29.0 |

Table 1: Dataset statistics about the number of documents for the train (#Train), development (#Dev) and test (#Test), the number (#Total) and the multi-event ratio (MER) of all documents.

| Precision | Recall | F1 | MER (%) |
|-----------|--------|------|---------|
| 98.8 | 89.7 | 94.0 | 31.0 |

Table 2: The quality of the DS-based event labeling evaluated on 100 manually annotated documents (randomly select 20 for each event type).

Doc2EDAG: 实验结果



| Model | EF | | | ER | | | EU | | | EO | | | EP | | |
|-----------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | P. | R. | F1 |
| DCFEE-O | 66.0 | 41.6 | 51.1 | 84.5 | 81.8 | 83.1 | 62.7 | 35.4 | 45.3 | 51.4 | 42.6 | 46.6 | 64.3 | 63.6 | 63.9 |
| DCFEE-M | 51.8 | 40.7 | 45.6 | 83.7 | 78.0 | 80.8 | 49.5 | 39.9 | 44.2 | 42.5 | 47.5 | 44.9 | 59.8 | 66.4 | 62.9 |
| GreedyDec | 79.5 | 46.8 | 58.9 | 83.3 | 74.9 | 78.9 | 68.7 | 40.8 | 51.2 | 69.7 | 40.6 | 51.3 | 85.7 | 48.7 | 62.1 |
| Doc2EDAG | 77.1 | 64.5 | 70.2 | 91.3 | 83.6 | 87.3 | 80.2 | 65.0 | 71.8 | 82.1 | 69.0 | 75.0 | 80.0 | 74.8 | 77.3 |

Table 3: Overall event-level precision (P.), recall (R.) and F1 scores evaluated on the test set.

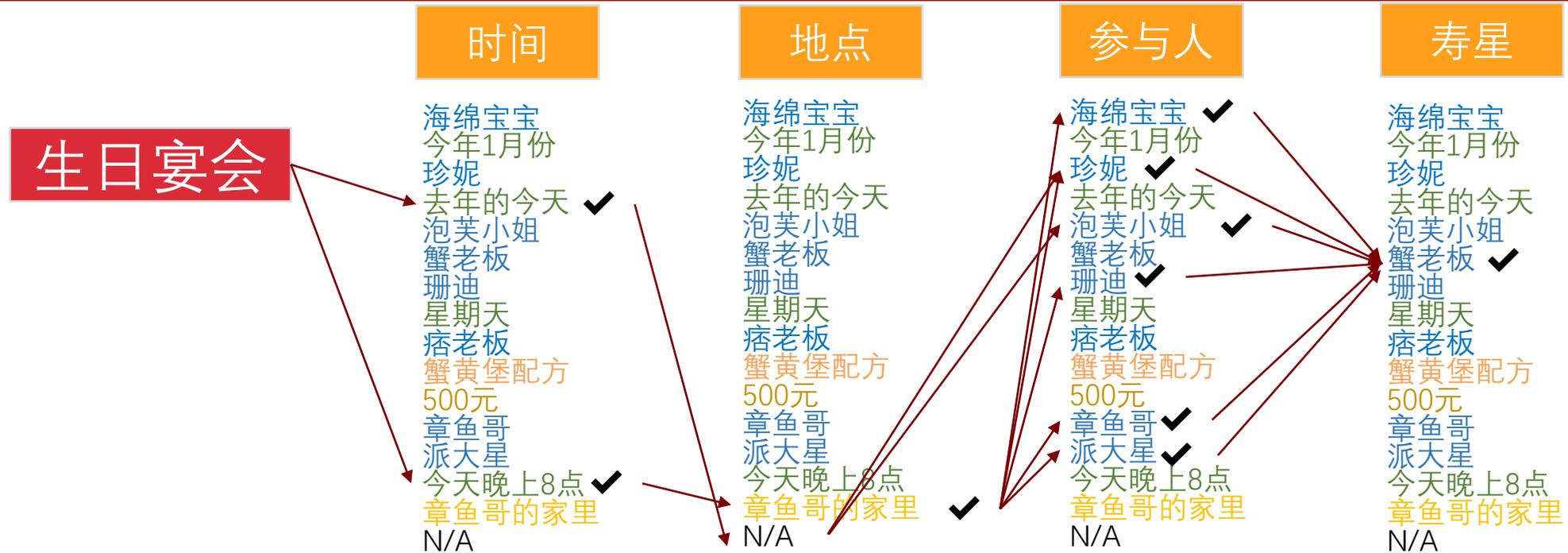
| Model | EF | | ER | | EU | | EO | | EP | | Avg. | | |
|-----------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | S. | M. | S. & M. |
| DCFEE-O | 56.0 | 46.5 | 86.7 | 54.1 | 48.5 | 41.2 | 47.7 | 45.2 | 68.4 | 61.1 | 61.5 | 49.6 | 58.0 |
| DCFEE-M | 48.4 | 43.1 | 83.8 | 53.4 | 48.1 | 39.6 | 47.1 | 42.0 | 67.0 | 60.6 | 58.9 | 47.7 | 55.7 |
| GreedyDec | 75.9 | 40.8 | 81.7 | 49.8 | 62.2 | 34.6 | 65.7 | 29.4 | 88.5 | 42.3 | 74.8 | 39.4 | 60.5 |
| Doc2EDAG | 80.0 | 61.3 | 89.4 | 68.4 | 77.4 | 64.6 | 79.4 | 69.5 | 85.5 | 72.5 | 82.3 | 67.3 | 76.3 |

Table 4: F1 scores for all event types and the averaged ones (Avg.) on single-event (S.) and multi-event (M.) sets.

- ER事件效果非常好
- 多实例中的表现不佳



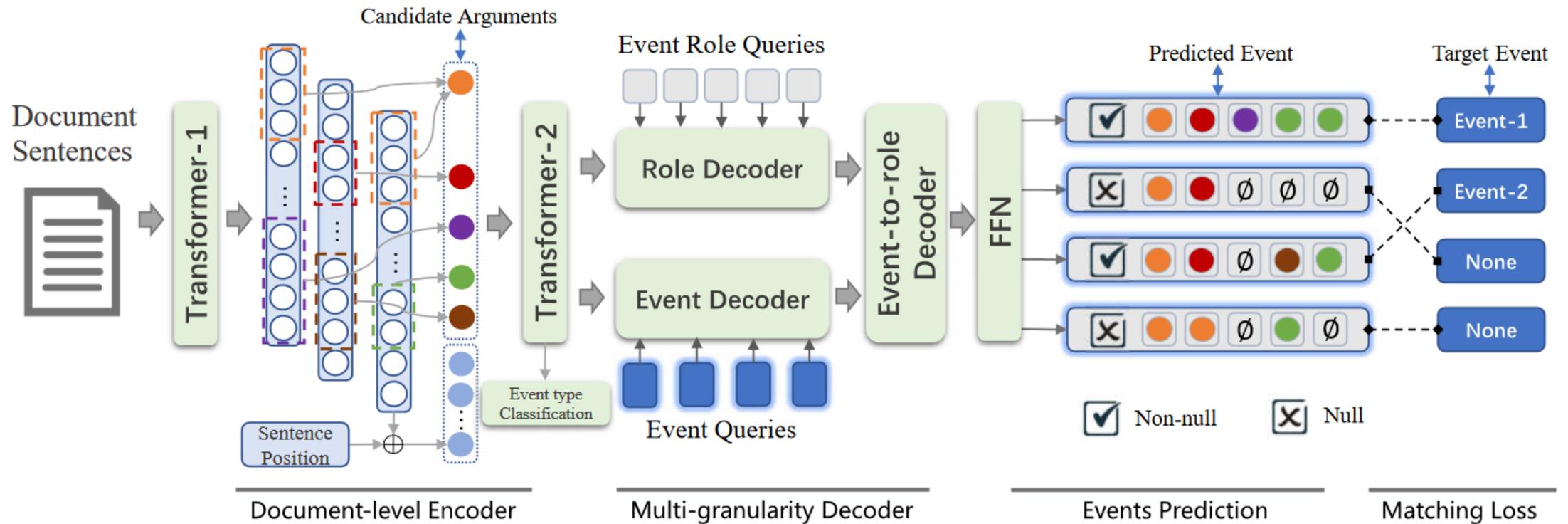
Doc2EDAG: 存在的不足



- 1、训练时动态生成EDAG，占用了很大的显存和内存空间
- 2、模型较大，训练时需要4-8张卡跑一星期，推理也非常慢
- 3、无法解决“参与人”这种一个要素角色对应多个要素的情况

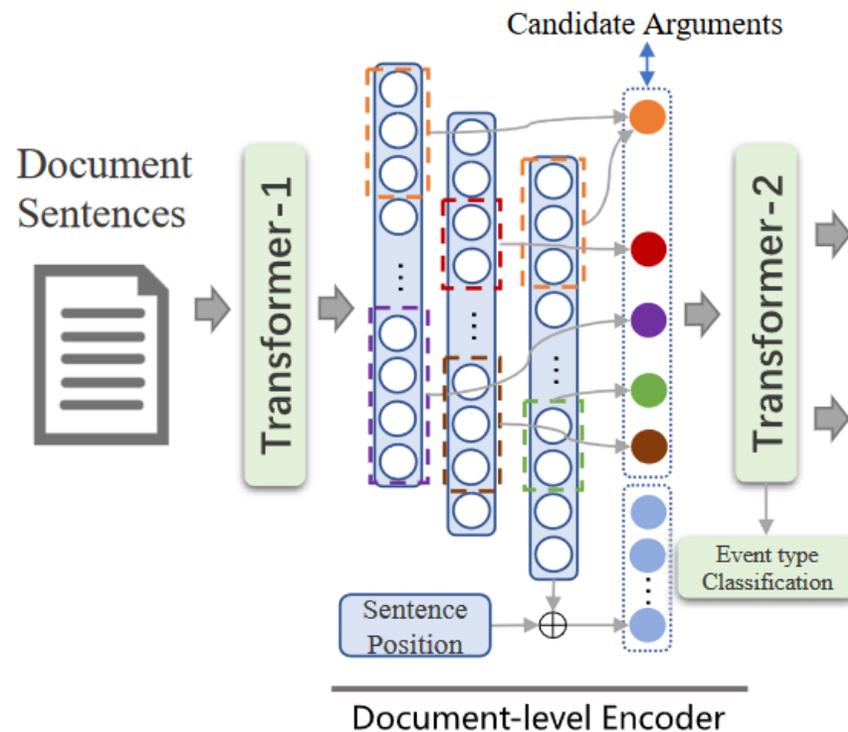
DE-PPN

- 假设存在 m 个事件实例，数据集中共有 n 种要素角色，抽取出 N'_a 个实体，则最终的表示结果可以用一个 $m \times n \times (N'_a + 1)$ 的张量表示，其中 +1 表示该位置的要素可能为 N/A



DE-PPN: 编码部分

- 假设存在 m 个事件实例，数据集中共有 n 种要素角色，抽取出 N'_a 个实体，则最终的结果可以用一个 $m \times n \times (N'_a + 1)$ 的张量表示，其中 +1 表示该位置的要素可能为 N/A



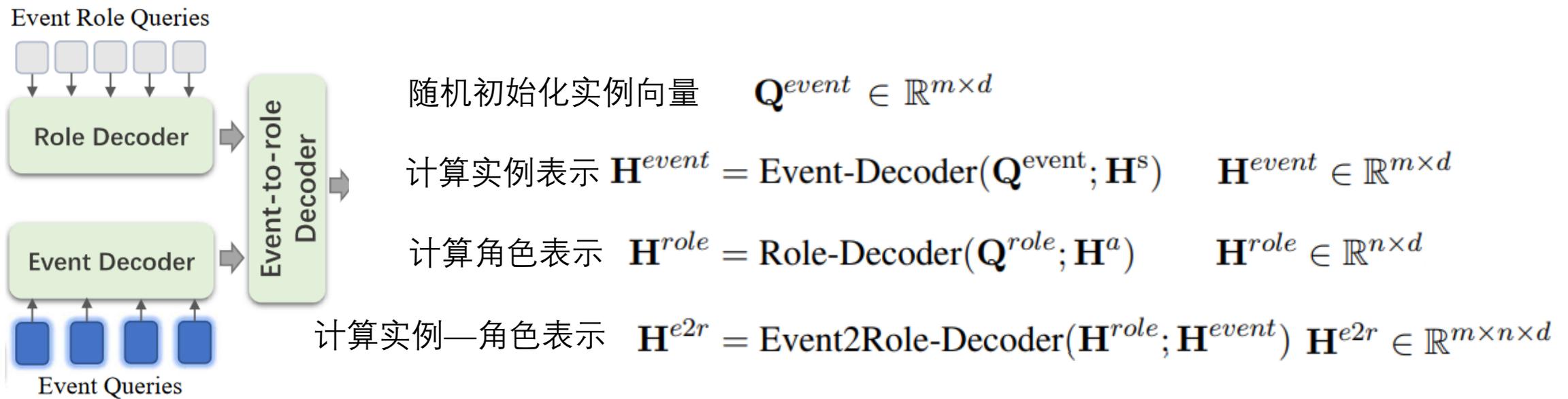
$$[\mathbf{H}^a; \mathbf{H}^s] = \text{Transformer-2}(\mathbf{c}_1^a \dots \mathbf{c}_{N'_a}^a; \mathbf{c}_1^s \dots \mathbf{c}_{N_s}^s)$$

$$\mathbf{H}^a \in \mathbb{R}^{N'_a \times d} \quad \mathbf{H}^s \in \mathbb{R}^{N_s \times d}$$

要素抽取、事件分类、要素表示增强部分和Doc2EDAG相同

DE-PPN: 表示求解

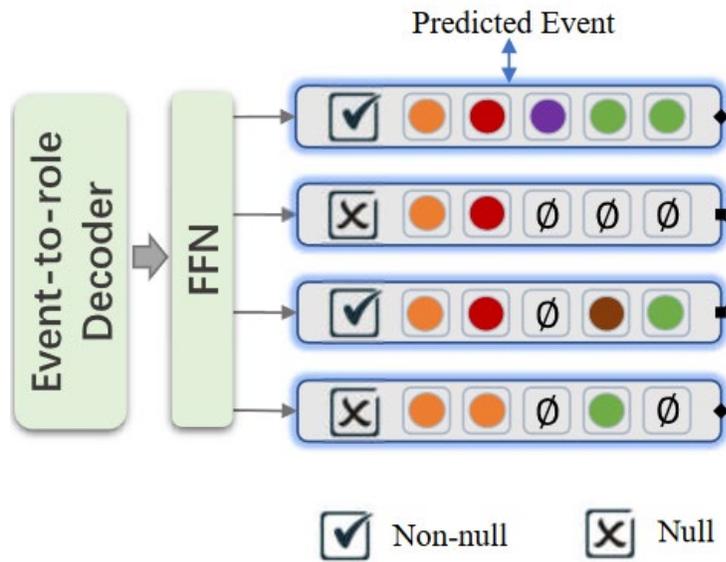
- 假设存在 m 个事件实例，数据集中共有 n 种要素角色，抽取出 N'_a 个实体，则最终的表示结果可以用一个 $m \times n \times (N'_a + 1)$ 的张量表示，其中 +1 表示该位置的要素可能为 N/A



Multi-granularity Decoder

DE-PPN: 解码部分

- 假设存在 m 个事件实例，数据集中共有 n 种要素角色，抽取出 N'_a 个实体，则最终的表示结果可以用一个 $m \times n \times (N'_a + 1)$ 的张量表示，其中 +1 表示该位置的要素可能为 N/A



判断生成的实例是否为空

$$\mathbf{p}^{event} = \text{softmax}(\mathbf{H}^{event} \mathbf{W}_e)$$

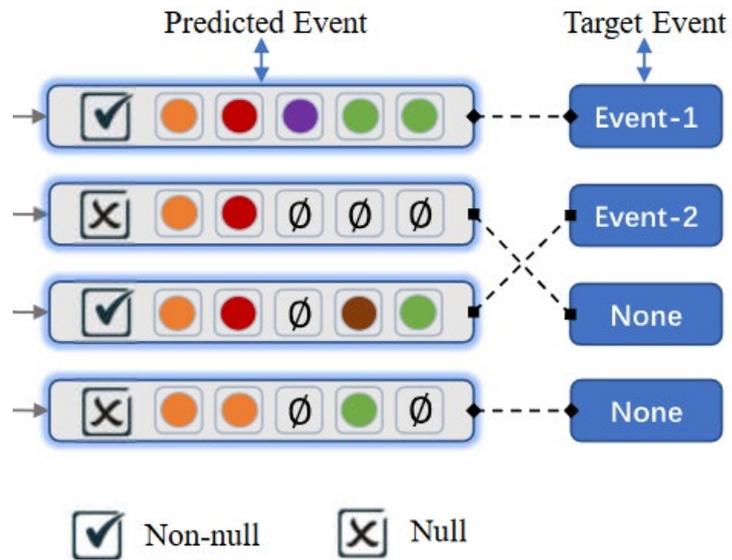
生成最终的实例结果

$$\mathbf{P}^{role} = \text{softmax}(\tanh(\mathbf{H}^{e2r} \mathbf{W}_1 + \mathbf{H}^a \mathbf{W}_2) \cdot \mathbf{v}_1)$$

$$\mathbf{P}^{role} \in \mathbb{R}^{m \times n \times (N'_a + 1)}$$

DE-PPN: 训练中的匹配问题

- 假设存在 m 个事件实例，数据集中共有 n 种要素角色，抽取出 N'_a 个实体，则最终的表示结果可以用一个 $m \times n \times (N'_a + 1)$ 的张量表示，其中 +1 表示该位置的要素可能为 N/A



$$\hat{\sigma} = \operatorname{argmax}_{\sigma \in \Pi(m)} \sum_i^m C_{\text{match}}(\hat{Y}_{\sigma(i)}, Y_i)$$

$$C_{\text{match}}(\hat{Y}_{\sigma(i)}, Y_i) = -\mathbb{1}_{\{\text{judge}_i \neq \emptyset\}} \sum_{j=1}^n P_{\sigma(i)}^j(r_i^j)$$

生成的 m 个事件实例需要与金标实例进行 loss 计算，但此时又无法得知和预测实例和答案之间的确切对应关系，因此引入线性匹配算法，并从金标中选择最可能对应的 m 个答案。



1. 抽取所有的要素

海绵宝宝早就等着这一天了。

今年1月份的时候，珍妮告诉他，自己的爸爸要过生日了，希望能给他办一个惊喜派对。

去年的今天，泡芙小姐给蟹老板办了一场惊喜派对，他非常开心。

但是珍妮手头没有足够的钱，希望海绵宝宝能够帮帮她。

海绵宝宝答应了她，并和好朋友珊迪说了这件事情。

珊迪和海绵宝宝在星期天找到了痞老板，并把蟹黄堡配方以500元的价格卖给了他。

这下他们有足够的钱啦~

珊迪说他们应该叫更多的朋友一起参加，比如章鱼哥和派大星。

大家商量了一下，决定今天晚上8点在章鱼哥的家里举行生日派对。

而此时的蟹老板还不知道等待他的是怎样的一个惊喜。

1. 抽取所有的要素

海绵宝宝早就等着这一天了。
今年1月份的时候，珍妮告诉他，自己的爸爸要过生日了，希望能给他办一个惊喜派对。
去年的今天，泡芙小姐给蟹老板办了一场惊喜派对，他非常开心。
但是珍妮手头没有足够的钱，希望海绵宝宝能够帮帮她。
海绵宝宝答应了她，并和好朋友珊迪说了这件事情。
珊迪和海绵宝宝在星期天找到了痞老板，并把蟹黄堡配方以500元的价格卖给了他。
这下他们有足够的钱啦~
珊迪说他们应该叫更多的朋友一起参加，比如章鱼哥和派大星。
大家商量了一下，决定今天晚上8点在章鱼哥的家里举行生日派对。
而此时的蟹老板还不知道等待他的是怎样的一个惊喜。

2. 事件分类

生日聚会、出售



DE-PPN: 例子

1. 抽取所有的要素

海绵宝宝早就等着这一天了。
 今年1月份的时候, 珍妮告诉他, 自己的爸爸要过生日了, 希望能给他办一个惊喜派对。
 去年的今天, 泡芙小姐给蟹老板办了一场惊喜派对, 他非常开心。
 但是珍妮手头没有足够的钱, 希望海绵宝宝能够帮帮她。
 海绵宝宝答应了她, 并和好朋友珊迪说了这件事情。
 珊迪和海绵宝宝在星期天找到了痞老板, 并把蟹黄堡配方以500元的价格卖给了他。
 这下他们有足够的钱啦~
 珊迪说他们应该叫更多的朋友一起参加, 比如章鱼哥和派大星。
 大家商量了一下, 决定今天晚上8点在章鱼哥的家里举行生日派对。
 而此时的蟹老板还不知道等待他的是怎样的一个惊喜。

2. 事件分类

生日聚会、出售

3. 事件生成

生日宴会#1

海绵宝宝
 今年1月份
 珍妮
 去年的今天
 泡芙小姐
 蟹老板
 珊迪
 星期天
 痞老板
 蟹黄堡配方
 500元
 章鱼哥
 派大星
 今天晚上8点
 章鱼哥的家里
 N/A

| 时间 | 地点 | 参与人 | 寿星 |
|----|----|-----|----|
| ✓ | | × | ✓ |
| | ✓ | | |

注: 这里用了
 argmax, 因此和
 Doc2EDAG一样, 一
 个角色不能对应多个
 要素



DE-PPN: 例子

1. 抽取所有的要素

海绵宝宝早就等着这一天了。
 今年1月份的时候, 珍妮告诉他, 自己的爸爸要过生日了, 希望能给他办一个惊喜派对。
 去年的今天, 泡芙小姐给蟹老板办了一场惊喜派对, 他非常开心。
 但是珍妮手头没有足够的钱, 希望海绵宝宝能够帮帮她。
 海绵宝宝答应了她, 并和好朋友珊迪说了这件事情。
 珊迪和海绵宝宝在星期天找到了痞老板, 并把蟹黄堡配方以500元的价格卖给了他。
 这下他们有足够的钱啦~
 珊迪说他们应该叫更多的朋友一起参加, 比如章鱼哥和派大星。
 大家商量了一下, 决定今天晚上8点在章鱼哥的家里举行生日派对。
 而此时的蟹老板还不知道等待他的是怎样的一个惊喜。

2. 事件分类

生日聚会、出售

3. 事件生成

生日宴会#2

海绵宝宝
 今年1月份
 珍妮
 去年的今天
 泡芙小姐
 蟹老板
 珊迪
 星期天
 痞老板
 蟹黄堡配方
 500元
 章鱼哥
 派大星
 今天晚上8点
 章鱼哥的家里
 N/A

| 时间 | 地点 | 参与人 | 寿星 |
|----|----|-----|----|
| | | ✓ | |
| | | × | |
| | | × | ✓ |
| | | × | |
| ✓ | | × | |
| | ✓ | × | |

注: 这里用了
 argmax, 因此和
 Doc2EDAG一样, 一
 个角色不能对应多个
 要素

DE-PPN: 实验结果



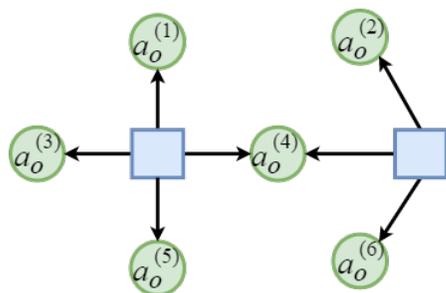
| Models | EF | | | ER | | | EU | | | EO | | | EP | | |
|-----------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | P | R | F1 |
| DCFEE-O | 66.0 | 41.6 | 51.1 | 84.5 | 81.8 | 83.1 | 62.7 | 35.4 | 45.3 | 51.4 | 42.6 | 46.6 | 64.3 | 63.6 | 63.9 |
| DCFEE-M | 51.8 | 40.7 | 45.6 | 83.7 | 78.0 | 80.8 | 49.5 | 39.9 | 44.2 | 42.5 | 47.5 | 44.9 | 59.8 | 66.4 | 62.9 |
| GreedyDec | 79.5 | 46.8 | 58.9 | 83.3 | 74.9 | 78.9 | 68.7 | 40.8 | 51.2 | 69.7 | 40.6 | 51.3 | 85.7 | 48.7 | 62.1 |
| Doc2EDAG | 77.1 | 64.5 | 70.2 | 91.3 | 83.6 | 87.3 | 80.2 | 65.0 | 71.8 | 82.1 | 69.0 | 75.0 | 80.0 | 74.8 | 77.3 |
| DE-PPN-1 | 77.8 | 55.8 | 64.9 | 75.6 | 76.4 | 76.0 | 76.4 | 63.7 | 69.4 | 77.1 | 54.3 | 63.7 | 85.5 | 43.0 | 57.2 |
| DE-PPN | 78.2 | 69.4 | 73.5 | 89.3 | 85.6 | 87.4 | 69.7 | 79.9 | 74.4 | 81.0 | 71.3 | 75.8 | 83.8 | 73.7 | 78.4 |

Table 1: Overall event-level precision (P), recall (R) and F1-score (F1) evaluated on the test set.

m=1

| Models | EF | | ER | | EU | | EO | | EP | | Avg. | | |
|-----------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | S. | M. | S.& M. |
| DCFEE-O | 56.0 | 46.5 | 86.7 | 54.1 | 48.5 | 41.2 | 47.7 | 45.2 | 68.4 | 61.1 | 61.5 | 49.6 | 58.0 |
| DCFEE-M | 48.4 | 43.1 | 83.8 | 53.4 | 48.1 | 39.6 | 47.1 | 42.0 | 67.0 | 60.0 | 58.9 | 47.7 | 55.7 |
| GreedyDec | 75.9 | 40.8 | 81.7 | 49.8 | 62.2 | 34.6 | 65.7 | 29.4 | 88.5 | 42.3 | 74.8 | 39.4 | 60.5 |
| Doc2EDAG | 80.0 | 61.3 | 89.4 | 68.4 | 77.4 | 64.6 | 79.4 | 69.5 | 85.5 | 72.5 | 82.3 | 67.3 | 76.3 |
| DE-PPN-1 | 82.4 | 46.3 | 78.3 | 53.9 | 82.2 | 45.6 | 78.1 | 39.3 | 82.8 | 38.5 | 80.7 | 44.7 | 66.2 |
| DE-PPN | 82.1 | 63.5 | 89.1 | 70.5 | 79.7 | 66.7 | 80.6 | 69.6 | 88.0 | 73.2 | 83.9 | 68.7 | 77.9 |

Table 2: F1-score for all event types and the averaged ones (Avg.) on single-event (S.) and multi-event (M.) sets.



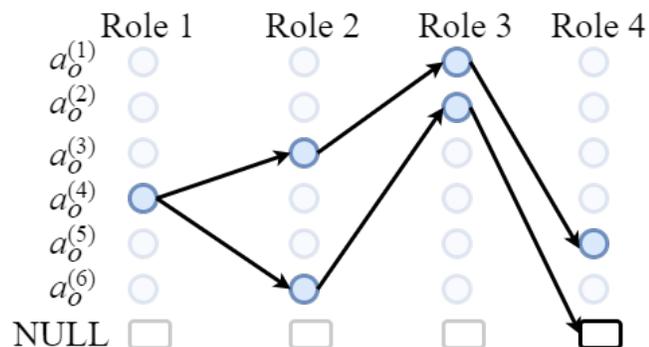
Annotated-Trigger-centered Trees

- 篇章事件抽取任务

- 触发词缺失或触发词被共用
- 长度很长，通常超过512个token

- 前人工作的痛点

- 训练的时候非常非常慢
 - 4-8 GPUs 跑一星期
- 需要巨大的内存和显存
 - 动态EDAG的历史空间

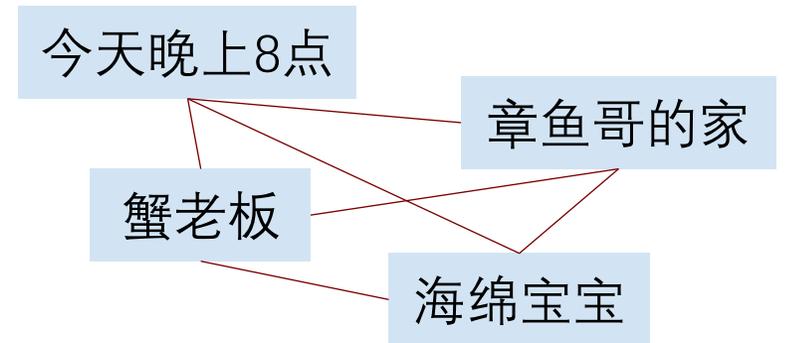
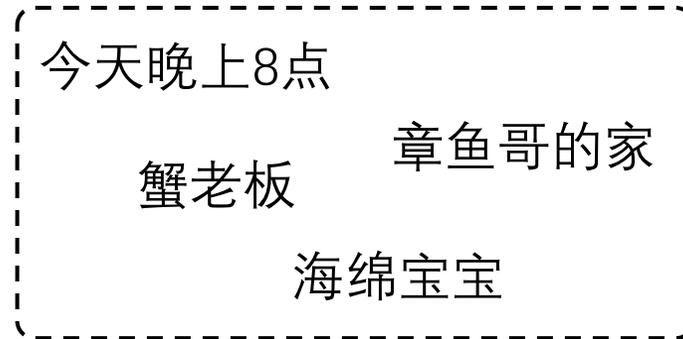


Directed Acyclic Graph
w/o Annotated Triggers

结合任务特点，针对以上痛点，提出了一种非自回归的 **快速** 和 **轻量化** 模型方案

PTPCG: 基于完全图的要素组合方案

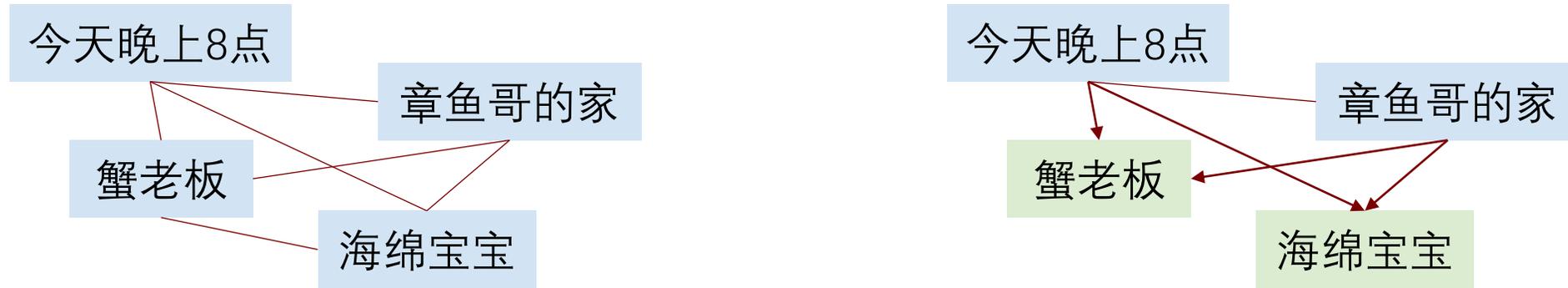
| 生日宴会 | |
|------|--------|
| 时间 | 今天晚上8点 |
| 地点 | 章鱼哥的家 |
| 寿星 | 蟹老板 |
| 参与人 | 海绵宝宝 |
| 主菜 | N/A |



假设

处于同一组合中的要素在语义空间中相近，因此它们之间应该互相双向连接，形成一张完全图。

PTPCG: 基于剪枝完全图的要素组合方案

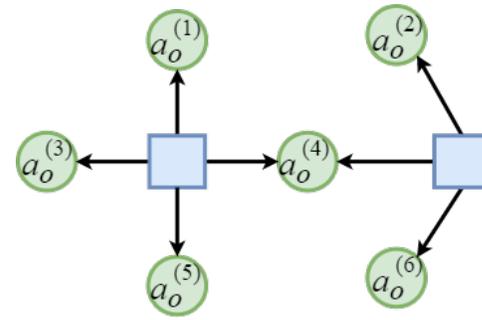
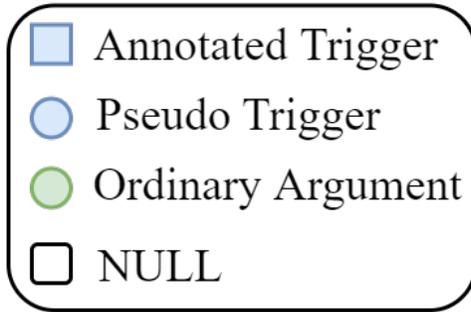


在实际应用中，不同要素的**重要性**不同。

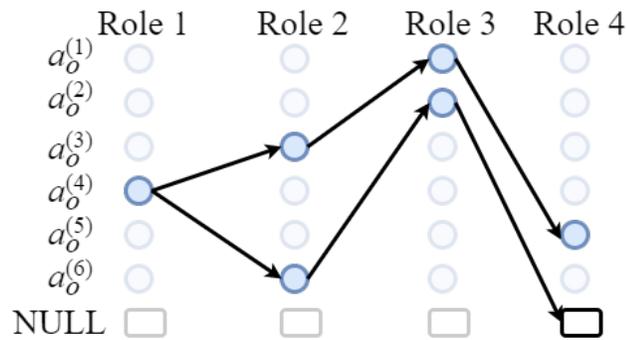
因此放松了上述假设的约束，只挑选最“重要”的要素作为“**伪触发词**”。

此时，无向（双向）完全图就被**剪枝**为了**有向图**。

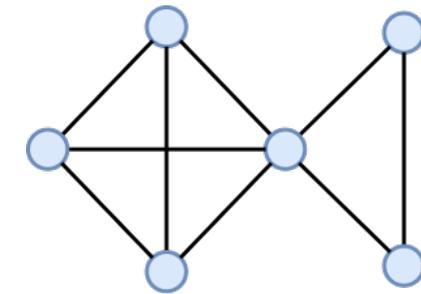
PTPCG: 不同要素组合方案之间的对比



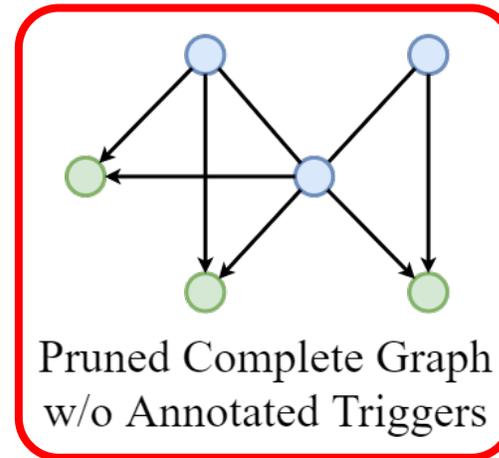
每个组合中只有一个伪触发词



Directed Acyclic Graph
w/o Annotated Triggers



Complete Graph
w/o Annotated Triggers

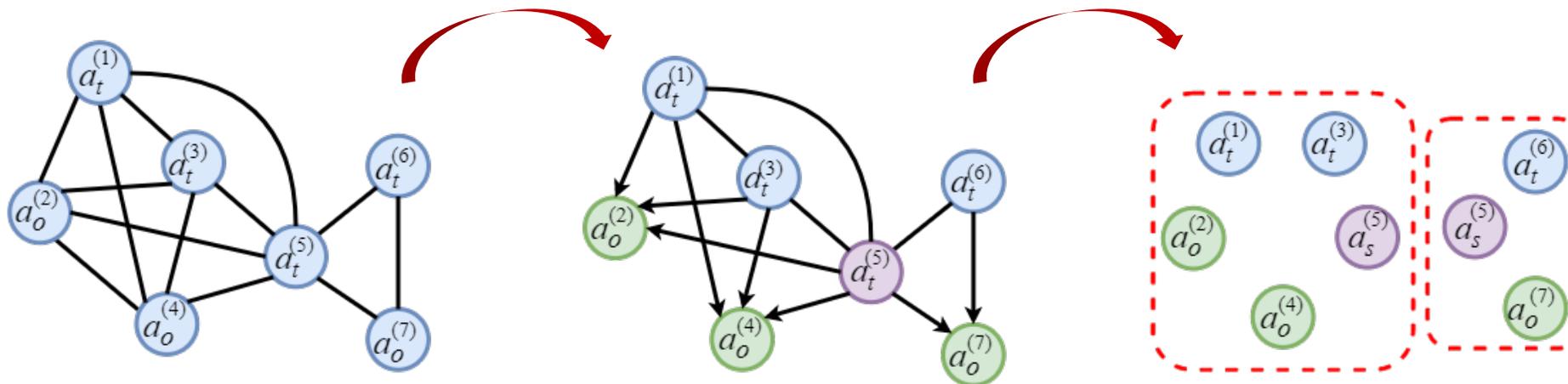


Pruned Complete Graph
w/o Annotated Triggers

剪枝

怎么挑选伪触发词

怎么解码成要素组合





PTPCG: 伪触发词选择

触发词在事件抽取中到底承担什么样的作用？

- **存在性 (Existence)** : 触发词用于指示事件实例
- **区分性 (Distinguishability)** :
 - 触发词在不同的事件实例中互不共享
 - 触发词可用于区分不同的事件实例

事件要素可以被认作伪触发词

对于事件 t_i , 伪触发词和一组事件要素角色的集合 \mathcal{R} 相关联
问题转变为: 怎样寻找合适的 \mathcal{R} ?



伪触发词选择

| | 🐼 | 🐼 | 🐼 |
|--------------------------|-----------------|------------|-----------------|
| Document \mathcal{D}_1 | 🍌 Plankton | 🍌 Plankton | 🍌 Squidward |
| | ★ Krabs | ★ Sandy | ★ Pearl |
| | 🦀 NULL | 🦀 NULL | 🦀 3,456,000 |
| | 👉 Dec. 16, 2016 | 👉 NULL | 👉 Nov. 16, 2016 |
| Document \mathcal{D}_2 | 🍌 Patrick | 🍌 Gary | 🍌 NULL |
| | ★ NULL | ★ Patrick | ★ SpongeBob |
| | 🦀 6,800 | 🦀 NULL | 🦀 6,800 |
| | 👉 Jan. 16, 1993 | 👉 NULL | 👉 Dec. 16, 2016 |

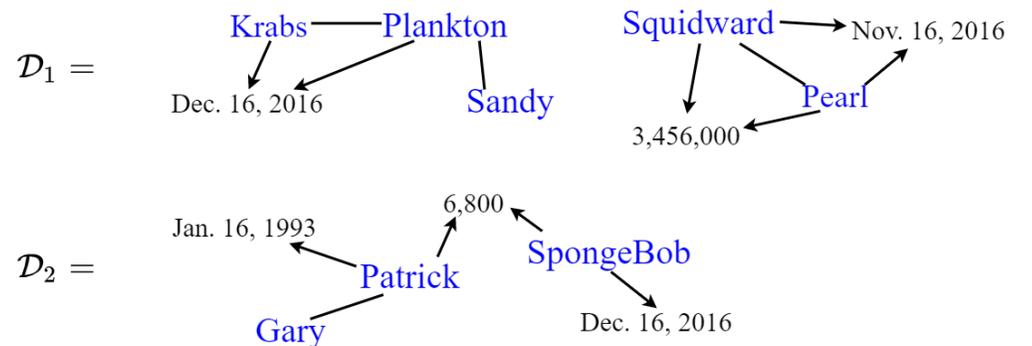
Existence(\mathcal{R}) = $\frac{\text{每篇文档中}\mathcal{R}\text{中至少有一个要素不为空的事件实例数量}}{\text{事件实例的总数}}$

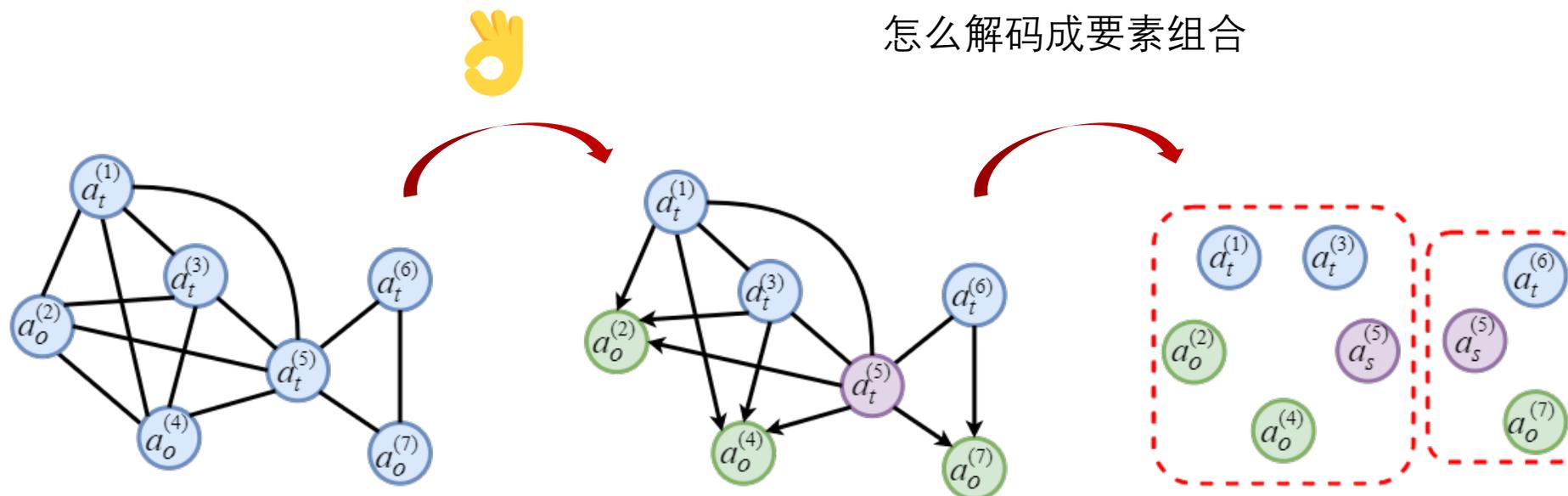
Distinguishability(\mathcal{R}) = $\frac{\text{每篇文档中}\mathcal{R}\text{中任一要素均不在其它实例中出现的实例数量}}{\text{事件实例的总数}}$

Importance(\mathcal{R}) = Existence(\mathcal{R}) × Distinguishability(\mathcal{R})

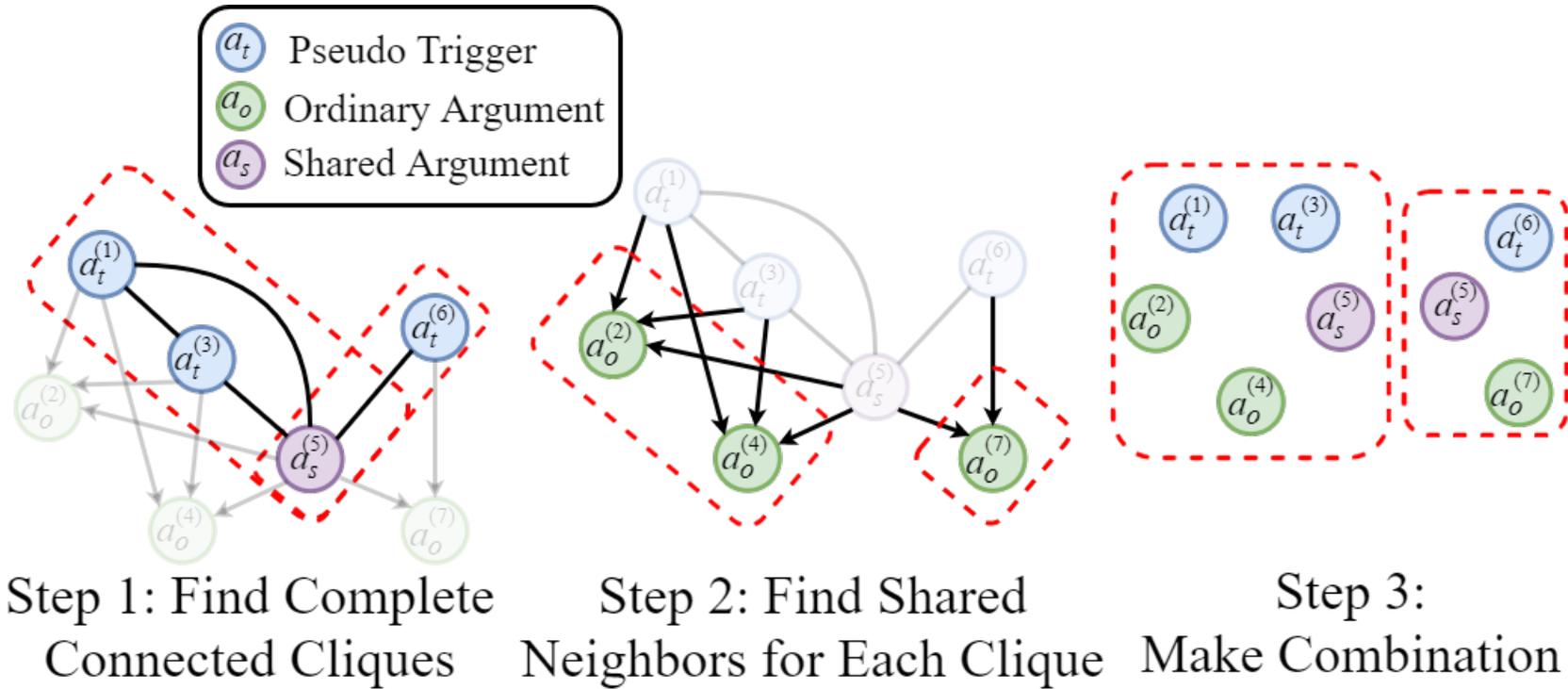
For $|\mathcal{R}| = 2$
 = 6 records

| Roles | Existence | Distinguishability | Importance | Sort ▼ |
|--------|-------------------------|-------------------------|--------------------------------------|------------|
| {🍌, ★} | $\frac{1+1+1+1+1+1}{6}$ | $\frac{1+1+1+1+1+1}{6}$ | $1 \times 1 = 1$ | → {🍌, ★} ✓ |
| {🍌, 🦀} | $\frac{1+1+1+1+1+1}{6}$ | $\frac{0+0+1+1+1+1}{6}$ | $1 \times \frac{4}{6} = \frac{2}{3}$ | → {🍌, 🦀} |
| {🍌, 👉} | $\frac{1+1+1+1+1+1}{6}$ | $\frac{1+1+1+1+1+1}{6}$ | $1 \times 1 = 1$ | → {🍌, 👉} |
| {★, 🦀} | $\frac{1+1+1+1+1+1}{6}$ | $\frac{1+1+1+1+1+1}{6}$ | $1 \times 1 = 1$ | → {★, 🦀} |
| {★, 👉} | $\frac{1+1+1+1+1+1}{6}$ | $\frac{1+1+1+1+1+1}{6}$ | $1 \times 1 = 1$ | → {★, 👉} |
| {🦀, 👉} | $\frac{1+0+1+1+0+1}{6}$ | $\frac{1+1+1+1+1+1}{6}$ | $\frac{4}{6} \times 1 = \frac{2}{3}$ | → {🦀, 👉} |

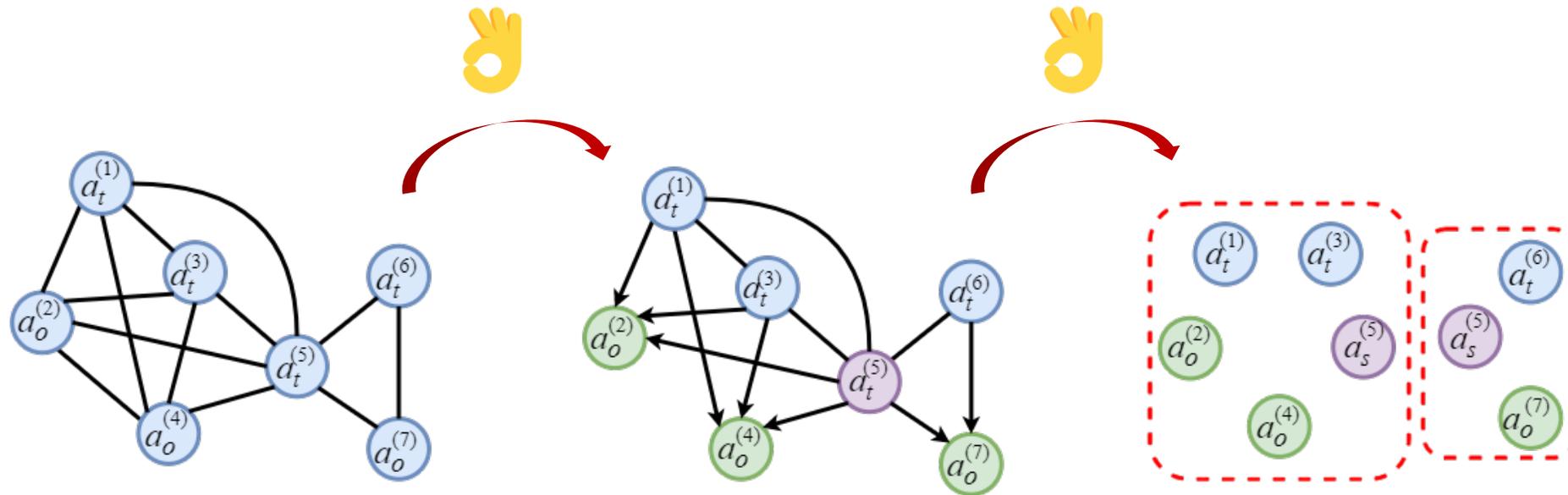




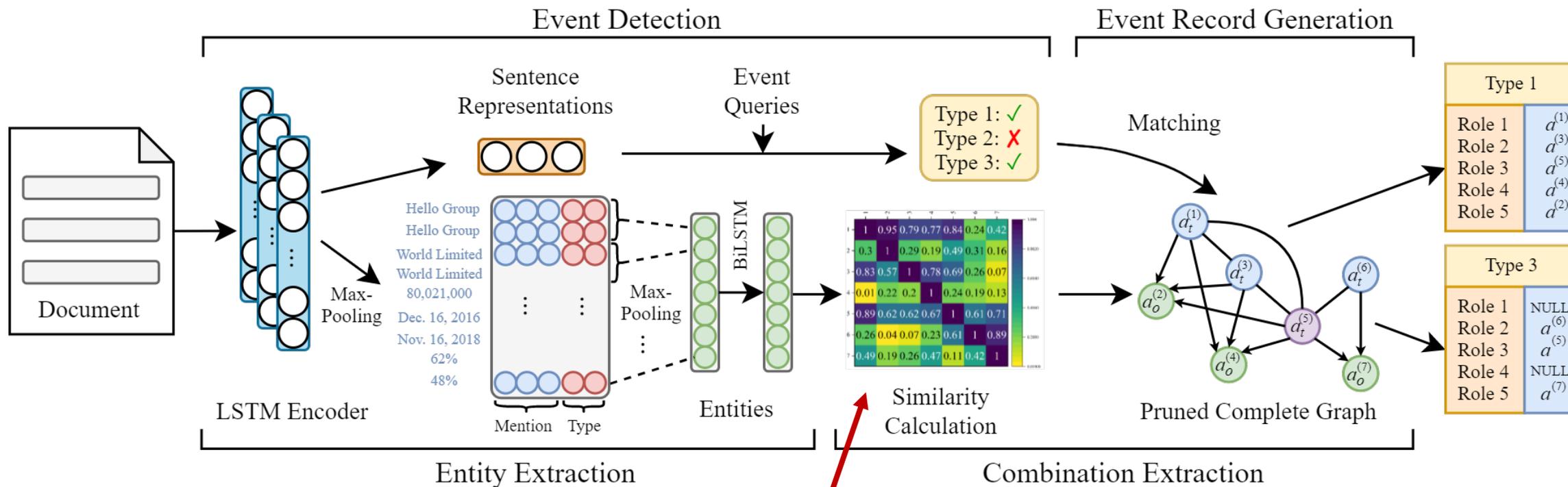
PTPCG: 组合解码



Bron-Kerbosch算法



PTPCG: 模型



dot-scaled attention

PTPCG: 实验结果



EDAG

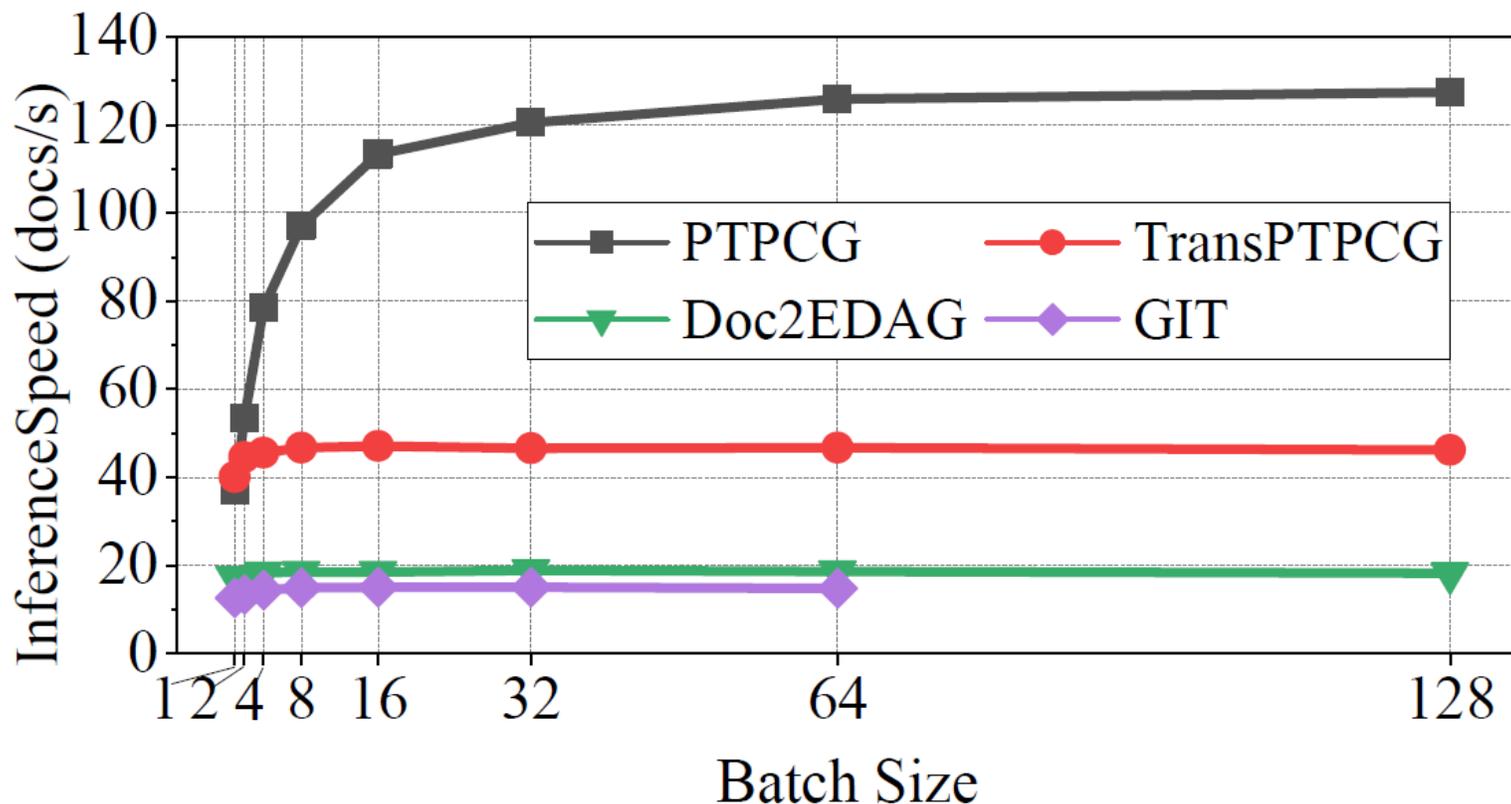
| Model | #Params (w/o Emb) | GPU Hours | ChFinAnn-Single | | | ChFinAnn-All | | | DuEE-fin w/o Tgg | | | DuEE-fin w/ Tgg | | |
|------------------------|----------------------|--------------|-----------------|-------------|-------------|--------------|-------------|-------------|------------------|-------------|-------------|-----------------|-------------|-------------|
| | | | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 |
| DCFEE-O* | 32M (16M) | 192.0 | 73.2 | 71.6 | 72.4 | 69.7 | 57.8 | 63.2 | 56.2 | 48.2 | 51.9 | 51.9 | 49.6 | 50.7 |
| DCFEE-M* | 32M (16M) | 192.0 | 64.9 | 71.7 | 68.1 | 60.1 | 61.3 | 60.7 | 38.7 | 52.3 | 44.5 | 37.3 | 48.6 | 42.2 |
| GreedyDec* | 64M (48M) | 604.8 | 83.9 | 77.3 | 80.4 | 81.9 | 51.2 | 63.0 | 59.6 | 41.8 | 49.1 | 59.0 | 42.1 | 49.2 |
| Doc2EDAG* | 64M (48M) | 604.8 | 83.2 | 89.3 | 86.2 | 81.1 | 77.0 | 79.0 | 66.7 | 50.0 | 57.2 | 67.1 | 51.3 | 58.1 |
| GIT* | 97M (81M) | 633.6 | 85.0 | 88.7 | 86.8 | 82.4 | 77.6 | 79.9 | 68.2 | 43.4 | 53.1 | 70.3 | 46.0 | 55.6 |
| PTPCG _{R =1} | 32M (16M) | 24.0 | 86.3 | 90.1 | 88.2 | 83.7 | 75.4 | 79.4 | 64.5 | 56.6 | 60.3 | 63.6 | 53.4 | 58.1 |

- **效果好**: PTPCG 和基于 EDAG的方案相比, 整体F1值持平, 在单事件单实例中的效果更佳
- **轻量**: PTPCG 只有 GIT **19.8%** 的参数
- **非常非常快**: *Blazing fast!* 单卡训练24小时即可, 只需要 GIT **3.9%** 的卡时

PTPCG 一个模型即可节省 **¥3658** ! (From ¥3801.6 to **¥144** !) *

*¥6/hour

PTPCG: 推理速度



比GIT 推理速度快 **8.5x**

- 目前还没有看到比 PTPCG 推理速度更快的方案
- 大 batch size 速度更佳
- 相较于基于 EDAG 的方案, 可拓展性更高



PTPCG: 伪触发词作为补充

| Model | Tgg | \mathcal{R} | Impt. | Dev | | | Online Test | | |
|----------|-----|---------------|-------|-------------|-------------|-------------|-------------|-------------|-------------|
| | | | | P | R | F1 | P | R | F1 |
| Doc2EDAG | × | - | - | 70.8 | 55.3 | 62.1 | 66.7 | 50.0 | 57.2 |
| | ✓ | - | - | 73.7 | 59.8 | 66.0 | 67.1 | 51.3 | 58.1 |
| GIT | × | - | - | 72.4 | 58.4 | 64.7 | 68.2 | 43.4 | 53.1 |
| | ✓ | - | - | 75.4 | 61.4 | 67.7 | 70.3 | 46.0 | 55.6 |
| PTPCG | ✓ | 0 | 62.9 | 73.5 | 59.4 | 65.7 | 67.0 | 50.1 | 57.3 |
| | ✓ | 1 | 93.7 | 68.8 | 64.2 | 66.4 | 62.0 | 54.8 | 58.1 |
| | ✓ | 2 | 97.1 | 64.7 | 64.9 | 64.8 | 59.1 | 56.5 | 57.8 |
| | × | 1 | 83.8 | 71.0 | 61.7 | 66.0 | 66.7 | 54.6 | 60.0 |
| | × | 2 | 94.3 | 63.8 | 64.8 | 64.3 | 60.2 | 58.4 | 59.3 |
| | × | 3 | 97.2 | 56.7 | 64.3 | 60.3 | 52.6 | 58.9 | 55.6 |

②

③

- ① DuEE-Fin 数据集拥有触发词标注，但没提供位置信息，存在共享的触发词
- ② 加上伪触发词之后，效果还能提升
- ③ 只有伪触发词的效果要比只有金标触发词的效果还好

We have a demo !



Visualisation on DocEE

Example Model

Text

Predicted Types

<https://github.com/Spico197/DocEE>

<http://hlt.suda.edu.cn/docee>

- 事件要素组合方案
 - 基于中心—卫星句: DCFEE
 - 基于自回归EDAG: Doc2EDAG
 - 基于Set: DE-PPN
 - 基于剪枝完全图: PTPCG
- 篇章事件抽取中, 由于标注较为困难, 因此常常存在触发词缺失或触发词被共用的情况
- 在触发词缺失的情况下, 传统基于触发词的要素组合方法不再适用
- 目前面向多实例的要素组合方案仍有很大的提升空间, 且各种事件类别之间尚存在较大的差异



- Cuiyun Han, Jinchuan Zhang, Xinyu Li, Guojin Xu, Weihua Peng and Zengfeng Zeng. 2022. DuEE-Fin: A Large-Scale Dataset for Document-Level Event Extraction. CCF International Conference on Natural Language Processing and Chinese Computing (NLPCC-22). pp 172–183
- Hang Yang, Yubo Chen, Kang Liu, Yang Xiao, and Jun Zhao. 2018. DCFEE: A Document-level Chinese Financial Event Extraction System based on Automatically Labeled Training Data. In Proceedings of ACL 2018, System Demonstrations, pages 50–55, Melbourne, Australia. Association for Computational Linguistics.
- Shun Zheng, Wei Cao, Wei Xu, and Jiang Bian. 2020. Doc2EDAG: An end-to-end document-level framework for Chinese financial event extraction. EMNLP-IJCNLP 2019 - 2019 Conference on Empirical Methods in Natural Language Processing and 9th International Joint Conference on Natural Language Processing, Proceedings of the Conference:337–346.
- Hang Yang, Dianbo Sui, Yubo Chen, Kang Liu, Jun Zhao, and Taifeng Wang. 2021. Document-level Event Extraction via Parallel Prediction Networks. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 6298–6308, Online. Association for Computational Linguistics.
- Tong Zhu, Xiaoye Qu, Wenliang Chen, Zhefeng Wang, Baoxing Huai, Nicholas Yuan and Min Zhang. 2022. Efficient Document-level Event Extraction via Pseudo-Trigger-aware Pruned Complete Graph. In Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence (IJCAI-22) Main Track. Pages 4552-4558. <https://doi.org/10.24963/ijcai.2022/632>

Thanks

Q&A

特别鸣谢故事主人公: 🍪 📁 🦋 🍷 📱 📧 🌸 ☆

朱桐

tzhu1997@outlook.com

<https://github.com/Spico197/DocEE>