# 第八章 知识图谱

戴洪良

计算机科学与技术学院/人工智能学院

hongldai@nuaa.edu.cn

# 从文本中抽取知识

- 命名实体识别
- **细粒度实体分类**
- 关系抽取
- 事件抽取

# 细粒度实体分类（Fine-grained Entity Typing, FET）

3月20日消息，微软旗下语音识别子公司Nuance今日发布一款 AI 临床笔记软件，命名为DAX Express，主要面向医护人员。

**目标输出**

| 实体提及[1] | 位置 | 类别 |
|---|---|---|
| 3月20日 | (0, 3) | /日期 |
| 微软 | (8, 9) | /机构, /机构/公司, /机构/公司/科技公司 |
| Nuance | (20, 25) | /机构, /机构/公司, /机构/公司/科技公司 |
| DAX Express | (46, 56) | /产品, /产品/软件 |

- 给定一个细粒度实体类别体系，对文本中已识别出的实体进行分类
- 先用NER先识别出实体位置，再用FET对它们分类
- FET一般是多标签分类，即一个样本可以有多个类别标签

# FET的应用

- 除用于知识图谱构建，还可应用于：
  - 关系抽取（relation extraction）
    - 头尾实体的类别可作为确定它们关系的辅助信息
  - 实体链接（entity linking）
    - 根据实体类别辅助确定正确的实体
  - 指代消解（coreference resolution）
    - 指代消解目标：识别文本中哪些词指代了相同的实体
    - 同实体类别的词或短语更可能指代了同一个实体
  - 等

Onoe, Yasumasa, and Greg Durrett. "Interpretable Entity Representations through Large-Scale Typing." Findings of EMNLP. 2020.

# FET的难点

- 实体类别数多
  - 有些类别可能不易区分
  - 需要的训练数据多

| Dataset | FIGER | OntoNotes | BBN | UltraFine | CFET |
|---------|-------|-----------|-----|-----------|------|
| #Types  | 117   | 89        | 47  | 10K       | 7k   |

- 预测的类别应与上下文相关，或需从上下文推断出来

Rogers, the UW's leading scorer, will be a game-time decision.

/organization, /organization/sports_team

# FET数据集

- 普通细粒度实体分类
  - 使用人工设计的实体类别体系，一般组织成层次结构
  - 类别例子：/person, /person/politician, /person/actor, /location/city
  - 数据集：FIGER、OntoNotes、BBN、Few-NERD等

- 极细粒度实体分类（Ultra-fine entity typing, UFET）
  - 直接使用普通的词或短语作为类别标签
  - 类别例子：person, politician, actor, city, victim, criminal, company
  - 数据集：Ultrafine、CFET等
    - 其中CFET为中文数据集

# FET数据集

- 普通FET数据集统计信息

| Dataset | FIGER | OntoNotes | BBN |
|---|---|---|---|
| #documents | 18 | 76 | 459 |
| #mentions | 563 | 9,604 | 13,766 |
| #types | 113 | 89 | 47 |
| type hierarchy depth | 2 | 3 | 2 |

**文档数（#documents）和样本数（#mentions）只统计人工标注的测试集**

- 其中提出FIGER数据集的论文是第一个系统化地重点研究细粒度实体分类的工作，该数据集标注质量也较高，但人工标注的用于测试的样本数少

- BBN的测试集中，每个样本只标注单条类别路径的标签
  - 如/LOCATION, /LOCATION/REGION，但不会给一个样本同时标/LOCATION, /ORGANIZATION
- BBN数据集没有对/PERSON类别进行细分

# FET数据集

- OntoNotes (Gillick et al, 2014) 数据集的类别体系



| PERSON | LOCATION | ORGANIZATION | OTHER | |
|---|---|---|---|---|
| **artist** | **structure** | **company** | **art** | **language** |
| actor | airport | broadcast | broadcast | programming |
| author | government | news | film | language |
| director | hospital | **education** | music | **living thing** |
| music | hotel | **government** | stage | animal |
| **education** | restaurant | **military** | writing | **product** |
| student | sports facility | **music** | **event** | camera |
| teacher | theatre | **political party** | accident | car |
| **athlete** | **geography** | **sports league** | election | computer |
| **business** | body of water | **sports team** | holiday | mobile phone |
| **coach** | island | **stock exchange** | natural disaster | software |
| **doctor** | mountain | **transit** | protest | weapon |
| **legal** | **transit** | | sports event | **food** |
| **military** | bridge | | violent conflict | **heritage** |
| **political figure** | railway | | **health** | **internet** |
| **religious leader** | road | | malady | **legal** |
| **title** | **celestial** | | treatment | **religion** |
| | **city** | | **award** | **scientific** |
| | **country** | | **body part** | **sports & leisure** |
| | **park** | | **currency** | **supernatural** |

Gillick, D., Lazic, N., Ganchev, K., Kirchner, J., & Huynh, D. (2014). Context-dependent fine-grained entity type tagging. *arXiv preprint arXiv:1412.1820.*

8

# FET数据集

- 普通FET数据样例

| Mention Example | Dataset | Labels |
|---|---|---|
| ... **Edward McClain** has been sitting at his spot outside the grocery store for years ... | FIGER (GOLD) | /person |
| **Silicon Valley** heaved a sigh of relief yesterday. | BBN | /LOCATION/REGION /LOCATION |
| The new **beer**, introduced this week at a liquor industry convention, is imported ... | BBN | /SUBSTANCE/FOOD /SUBSTANCE |
| ... modestly compares its "hydraulic active suspension" to a **cheetah** ... | BBN | /ANIMAL |
| ... but **Valley Federal** had said it expected to post a modest pretax gain and to ... | OntoNotes | /organization/company /organization |
| ... but Valley Federal had said **it** expected to post a modest pretax gain and to ... | OntoNotes | /organization/company /organization |
| ... and **at least 500 civilians hiding inside** were killed, more than half of whom ... | OntoNotes | /person |
| ... and to promote **development** in other areas of the two countries. | OntoNotes | /other |

# 极细粒度实体分类

- 不人工设计类别体系，直接使用词或短语作为类别标签
  - 类别例子：person, politician, actor, city, victim, patient, company


- 相比普通FET，类别更丰富，覆盖更广
  - Ultrafine和CFET分别含约10k和7k实体类别

# FET数据集

- UFET数据样例

| 数据集 | 样本 | 标签 |
|---|---|---|
| Ultrafine | In 1988 , **Pitt** had his first starring role, in Dark Side Of The Sun, where he played a young American taken by his family to the ... | performer, adult, man, male, entertainer, professional, person, actor |
| Ultrafine | ... states that Paul of Tarsus, imprisoned and on trial claimed his right as a Roman citizen to be tried before Caesar, and the judicial process had to be suspended until **he** was brought to Rome. | citizen, criminal, person |
| CFET | 高尔基大街（现易名为**特维尔大街**）是莫斯科一条最主要的大街 | 街道, 旅游景点, 路, 大街, 街, 道路 |
| CFET | 我在**西堤**牛排上海虹口龙之梦店：同学小聚 哈哈 | 品牌, 地方, 餐馆, 位置 |

Ultrafine和CFET分别含约6000和4800个人工标注的样本，等分为train/dev/test

# 极细粒度实体分类

- 缺点：增加了应用的难度
  - 分类预测的效果变差，甚至人工标注数据时也可能标不全
    - 类别多，每个样本对应的正确类别标签也多
  - 基于分类结果执行某些操作的难度增加
  - 类别定义不清晰，有些类别词可能有歧义

# 自动构建训练数据

- FET类别数多，标注难度大，人工标注训练数据成本高
  - 应对方法：自动生成训练数据
  - 学术界对细粒度实体分类的研究目前大多基于自动生成的训练数据

- 三种主要的自动构建训练数据方法
  - 基于知识图谱的构建方法
  - 基于head word的方法
  - 基于预训练模型的方法

# 自动构建训练数据

- 基于知识图谱的构建方法
  1. 获得实体提及
     - 使用Wikipedia的内部超链接（anchor links），或使用NER工具标注
  2. 得到实体提及在知识图谱中的对应实体
     - 基于Wikipedia内部超链接的链接目标，或使用实体链接
  3. 基于知识图谱中实体的类别得到标签
     - 将知识图谱中的实体类别映射到所使用的标签体系

# 基于知识图谱的训练数据获取 – 例

**A piece of text from Wikipedia**

**Wikipedia page of Arnold Schwarzenegger**

increase. On November 17, Gov. Schwarzenegger signed Executive Order S-1-03, rescinding the vehicle license fee retroactive to October 1, 2003 when the fee increase went into effect. Analysts

**mention**

## Arnold Schwarzenegger

From Wikipedia, the free encyclopedia

**Arnold Alois Schwarzenegger** (/ˈʃvɑːrtsnɛɡər/;[1][a] German: [ˈaɐ̯nɔlt ˈʃvaɐ̯tsn̩ ˈʔɛɡɐ]; born July 30, 1947) is an Austrian-American actor, filmmaker, businessman, author, philanthropist, activist, politician, and former professional bodybuilder and powerlifter.[2] He served as the 38th Governor of California from 2003 to 2011.

Obtain Freebase types

**Target Types**:
*/person; /person/actor; /person/athlete; /person/politician; /person/author*

Map to target types

**Freebase Types**:
*people.person; film.actor; sports.pro_athlete; government.politician; book.author; …*

**weak labels for the mention**

这种训练数据也被叫做远程监督

# 基于知识图谱的训练数据获取

- 存在的问题：
  - 使用Wikipedia内部超链接或实体链接得到所指代的实体都可能出错
    - 从而根据实体从知识库中获得的类别标签也是错的
  - 得到的标签与上下文无关
  - 应用到UFET时召回率低
    - 在Ultrafine数据集上，该方法平均为每个样本获取的标签数少于2个，但人工标注的样本平均每个有5.4个标签

# 自动构建FET训练数据

- 基于head word生成

  - Head word: the central element in a phrase

如：

| phrase | head word |
| --- | --- |
| the 44th president of the United States | president |
| Nanjing University of Aeronautics and Astronautics | university |
| the man with her | man |
| a group of students | group |

**用于极细粒度实体分类**：如果head word是一个目标类别，则直接使用head word作为标签
**用于普通细粒度实体分类**：将head word映射到类别体系中的标签

# 基于head word生成类别标签

- 存在的问题：
  - 对很多实体提及不适用，如"Bob Dylan"、"Microsoft"
  - 应用到UFET时召回率低

# 自动构建训练数据

- 基于预训练语言模型的方法

An unlabeled mention in a sentence:

Highway 61 Revisited is the sixth studio album by Bob Dylan.

Highway 61 Revisited is the sixth studio album by [MASK] such as Bob Dylan.

Feed to BERT MLM

Most probable words for "[MASK]": artists, musicians, musician, songwriters, performers

*artist*, *musician*, *songwriter*, *performer*  ➡  作为标签

Dai, Hongliang, Yangqiu Song, and Haixun Wang. "Ultra-Fine Entity Typing with Weak Supervision from a Masked Language Model." Proceedings of ACL. 2021.

# 自动构建训练数据

- 基于预训练语言模型的方法

| Input | Top Words for [MASK] |
|---|---|
| In late 2015, [MASK] such as Leonardo DiCaprio starred in The Revenant. | actors, stars, actor, directors, filmmakers |
| At some clinics, they and some other [MASK] are told the doctors don't know how to deal with AIDS, and to go someplace else. | patients, people, doctors, kids, children |
| Finkelstein says he expects the company to "benefit from some of the disruption faced by our competitors and any other [MASK]." | company, business, companies, group, investors |

- 该方法为不同形式的实体提及（专有名词、代词、普通名词）都可生成类别标签
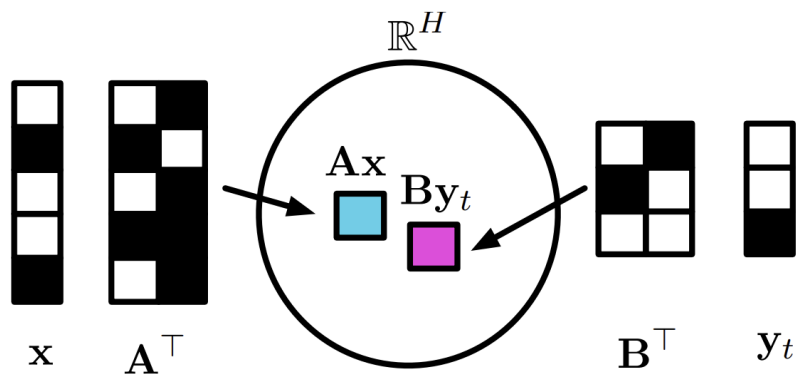- 可以生成像patient这样的依赖于上下文的标签

# 基于预训练模型的训练数据生成方法

- 优点
  - 可以得到与上下文相关的标签，且对实体提及的形式要求低
  - 可以补充基于知识图谱和基于head word的方法，提高生成标签的召回率

- 存在的问题
  - 得到的标签正确性相比基于知识图谱和head word的方法更低
  - 生成标签数目不易确定

| Pattern | F1 |
| --- | --- |
| $M$ and any other $H$ | 25.3 |
| $M$ and some other $H$ | 24.8 |
| $H$ such as $M$ | 20.7 |
| such $H$ as $M$ | 18.1 |
| $H$ including $M$ | 17.4 |
| $H$ especially $M$ | 11.5 |

# 细粒度实体分类方法

- 基于手工设计特征的方法

- (Yogatama et al., 2015)的方法

  - 将mention和类别嵌入到同一向量空间后求点积，得到mention是否属于该类别的分数



$$s(\mathbf{x}, \mathbf{y}_t; \mathbf{A}, \mathbf{B}) = f(\mathbf{x}, \mathbf{A}) \cdot g(\mathbf{y}_t, \mathbf{B}) = \mathbf{Ax} \cdot \mathbf{By}_t$$

**其中** $\boldsymbol{x}$：mention的特征向量；$\boldsymbol{y}_t$：类别标签$t$的one-hot encode向量
$\mathbf{A}$和$\mathbf{B}$：可训练参数矩阵

预测时，认为使分数$s$大于一个预设阈值的标签为正确标签

Yogatama, Dani, Dan Gillick, and Nevena Lazic. "Embedding methods for fine grained entity type classification." *Proceedings of ACL-IJCNLP*. 2015.

# 细粒度实体分类方法

- 基于手工设计特征的方法

- (Yogatama et al., 2015)使用的特征

| Feature | Description | Example |
|---|---|---|
| Head | The syntactic head of the mention phrase | "Obama" |
| Non-head | Each non-head word in the mention phrase | "Barack", "H." |
| Cluster | Word cluster id for the head word | "59" |
| Characters | Each character trigram in the mention head | ":ob", "oba", "bam", "ama", "ma:" |
| Shape | The word shape of the words in the mention phrase | "Aa A. Aa" |
| Role | Dependency label on the mention head | "subj" |
| Context | Words before and after the mention phrase | "B:who", "A:first" |
| Parent | The head's lexical parent in the dependency tree | "picked" |
| Topic | The most likely topic label for the document | "politics" |

# 细粒度实体分类方法

- 模型训练的loss函数（margin-based）

$$l(m_i, Y_i, \bar{Y}_i) = \sum_{y \in Y_i} \sum_{\bar{y} \in \bar{Y}_i} \max\left(0, \gamma - s(m_i, y) + s(m_i, \bar{y})\right)$$

$m_i$：第$i$个mention

$Y_i$：第$i$个mention的正确标签集合

$\bar{Y}_i$：第$i$个mention的不正确标签集合

$s(m_i, y)$：模型输出的$m_i$属于类别$y$的分数

# 细粒度实体分类方法

- 基于手工设计特征的方法

- (Yogatama et al., 2015)的实验效果

| Method | P | R | F1 |
|---|---|---|---|
| Ling and Weld (2012) | – | – | 69.30 |
| WSABIE | 81.85 | 63.75 | 71.68 |
| K-WSABIE | **82.23** | **64.55** | **72.35** |

FIGER 数据集上的效果（Micro-average Precision, Recall, F1）

- 指标计算：

$$P = \frac{\sum_{m \in M} |y_m \cap \hat{y}_m|}{\sum_{m \in M} |\hat{y}_m|} \qquad R = \frac{\sum_{m \in M} |y_m \cap \hat{y}_m|}{\sum_{m \in M} |y_m|}$$
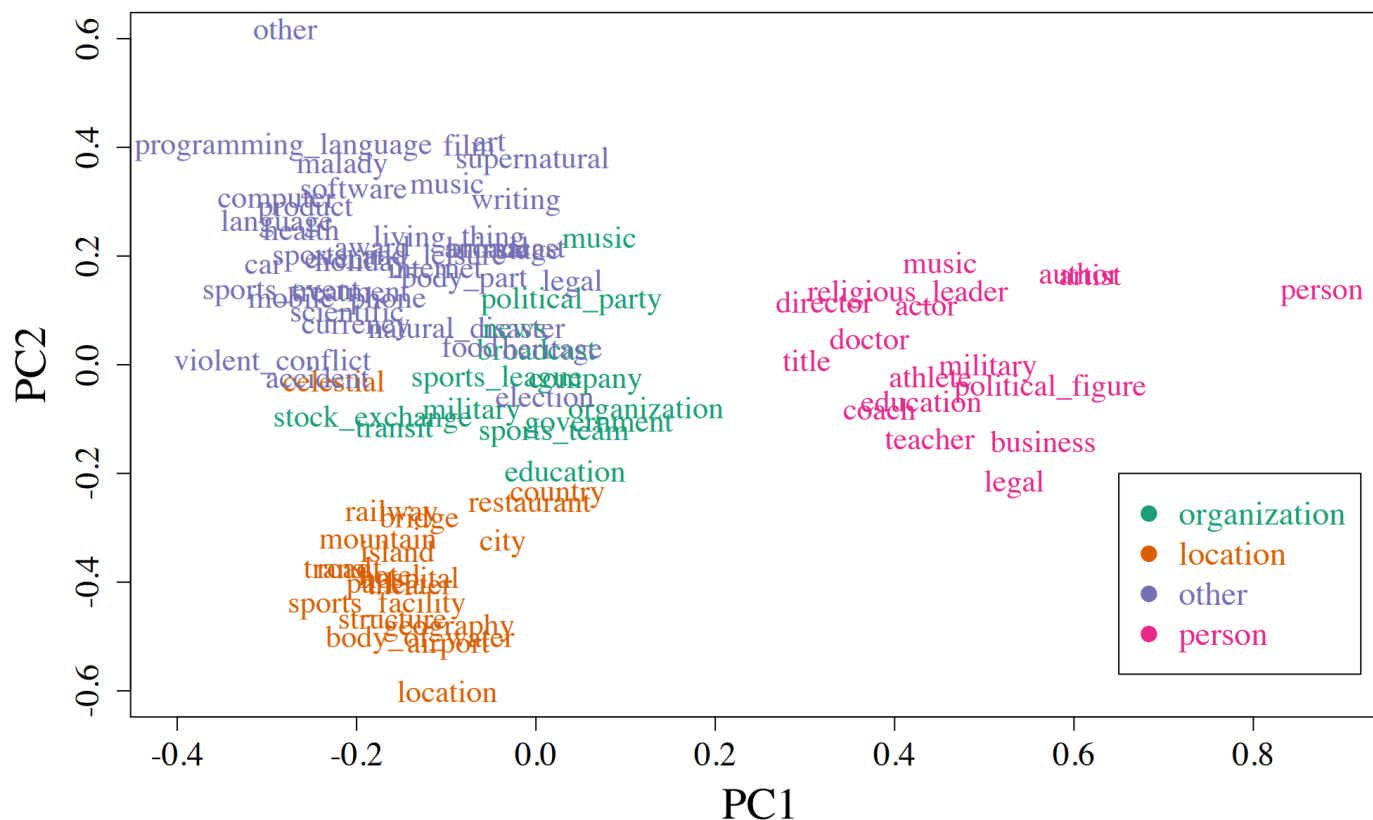
$$F1 = 2 * P * R / (P + R)$$

其中 $M$ 为测试集需预测类别的实体提及集合，$y_m$ 和 $\hat{y}_m$ 分别为实体提及 $m$ 的正确标签集合和预测标签集合

# 细粒度实体分类方法

- 基于手工设计特征的方法

- (Yogatama et al., 2015)的实验效果

将学习到的类别标签向量映射到2维空间：

属于同一大类的类别在向量空间中距离较近

# 基于神经网络的细粒度实体分类

- (Abhishek et al., 2017)的方法

Abhishek, A., Anand, A., & Awekar, A. (2017). Fine-Grained Entity Type Classification by Jointly Learning Representations and Label Embeddings. In *Proceedings of EACL* (Vol. 1, pp. 797-807).

# 基于神经网络的细粒度实体分类

- (Abhishek et al., 2017)的方法

**训练loss:**

**对标签属于同一类别路径的 (clean)：**

$$l(m_i, Y_i, \bar{Y}_i) = \sum_{y \in Y_i} \max(0, 1 - s(m_i, y)) + \sum_{y \in \bar{Y}_i} \max(0, 1 + s(m_i, y))$$

**对标签属于不同类别路径的 (noisy)：**

$$l(m_i, Y_i, \bar{Y}_i) = \max(0, 1 - s(m_i, y^*)) + \sum_{y \in \bar{Y}_i} \max(0, 1 + s(m_i, y))$$

其中　　$y^* = argmax_{y \in Y_i} s(m_i, y)$

对正标签，因为其中有些可能与上下文无关，只用其中得分最高的那个计算loss

# 细粒度实体分类方法

- (Abhishek et al., 2017)的实验效果

| Typing methods | Wiki/FIGER(GOLD) | | | OntoNotes | | | BBN | | |
|---|---|---|---|---|---|---|---|---|---|
| | Acc. | Ma-F1 | Mi-F1 | Acc. | Ma-F1 | Mi-F1 | Acc. | Ma-F1 | Mi-F1 |
| **FIGER**[*] (Ling and Weld, 2012) | 0.474 | 0.692 | 0.655 | 0.369 | 0.578 | 0.516 | 0.467 | 0.672 | 0.612 |
| **HYENA**[*] (Yosef et al., 2012) | 0.288 | 0.528 | 0.506 | 0.249 | 0.497 | 0.446 | 0.523 | 0.576 | 0.587 |
| **AFET-NoCo**[*] (Ren et al., 2016) | 0.526 | 0.693 | 0.654 | 0.486 | 0.652 | 0.594 | 0.655 | 0.711 | 0.716 |
| **AFET-CoH**[*] (Ren et al., 2016) | 0.433 | 0.583 | 0.551 | 0.521 | 0.680 | 0.609 | 0.657 | 0.703 | 0.712 |
| **AFET**[*] (Ren et al., 2016) | 0.533 | 0.693 | 0.664 | 0.551 | 0.711 | 0.647 | 0.670 | 0.727 | 0.735 |
| **AFET**[†‡] (Ren et al., 2016) | 0.509 | 0.689 | 0.653 | **0.553** | **0.712** | **0.646** | 0.683 | 0.744 | 0.747 |
| **Attentive**[†] (Shimaoka et al., 2016) | 0.581 | 0.780 | 0.744 | 0.473 | 0.655 | 0.586 | 0.484 | 0.732 | 0.724 |
| **our-AllC**[†] | **0.662** | 0.805 | 0.770 | 0.514 | 0.672 | 0.626 | 0.655 | 0.736 | 0.752 |
| **our-NoM**[†] | 0.646 | 0.808 | 0.768 | 0.521 | 0.683 | 0.626 | 0.615 | 0.742 | 0.755 |
| **our**[†] | 0.658 | **0.812** | **0.774** | 0.522 | 0.685 | 0.633 | 0.604 | 0.741 | 0.757 |
| **model level transfer-learning**[†] | - | - | - | 0.531 | 0.684 | 0.637 | 0.645 | 0.784 | **0.795** |
| **feature level transfer-learning**[†] | - | - | - | 0.471 | 0.689 | 0.635 | **0.733** | **0.791** | 0.792 |

**Our-AllC**:直接用未改进的loss
**Our-NoM**: 不用mention representation

**Acc**: 预测对的实体提及数/总实体提及数

**Ma-F1**计算:

$$P = \frac{1}{|M|} \sum_{m \in M} \frac{|y_m \cap \hat{y}_m|}{|\hat{y}_m|} \quad R = \frac{1}{|M|} \sum_{m \in M} \frac{|y_m \cap \hat{y}_m|}{|y_m|}$$

$$F1 = 2 * P * R/(P + R)$$

**Mi-F1**计算:

$$P = \frac{\sum_{m \in M} |y_m \cap \hat{y}_m|}{\sum_{m \in M} |\hat{y}_m|} \qquad R = \frac{\sum_{m \in M} |y_m \cap \hat{y}_m|}{\sum_{m \in M} |y_m|}$$

$$F1 = 2 * P * R/(P + R)$$

# 细粒度实体分类方法

- (Abhishek et al., 2017)的实验效果

| Typing methods | Wiki/FIGER(GOLD) | | | OntoNotes | | | BBN | | |
|---|---|---|---|---|---|---|---|---|---|
| | Acc. | Ma-F1 | Mi-F1 | Acc. | Ma-F1 | Mi-F1 | Acc. | Ma-F1 | Mi-F1 |
| **FIGER**[*] (Ling and Weld, 2012) | 0.474 | 0.692 | 0.655 | 0.369 | 0.578 | 0.516 | 0.467 | 0.672 | 0.612 |
| **HYENA**[*] (Yosef et al., 2012) | 0.288 | 0.528 | 0.506 | 0.249 | 0.497 | 0.446 | 0.523 | 0.576 | 0.587 |
| **AFET-NoCo**[*] (Ren et al., 2016) | 0.526 | 0.693 | 0.654 | 0.486 | 0.652 | 0.594 | 0.655 | 0.711 | 0.716 |
| **AFET-CoH**[*] (Ren et al., 2016) | 0.433 | 0.583 | 0.551 | 0.521 | 0.680 | 0.609 | 0.657 | 0.703 | 0.712 |
| **AFET**[*] (Ren et al., 2016) | 0.533 | 0.693 | 0.664 | 0.551 | 0.711 | 0.647 | 0.670 | 0.727 | 0.735 |
| **AFET**[†‡] (Ren et al., 2016) | 0.509 | 0.689 | 0.653 | **0.553** | **0.712** | **0.646** | 0.683 | 0.744 | 0.747 |
| **Attentive**[†] (Shimaoka et al., 2016) | 0.581 | 0.780 | 0.744 | 0.473 | 0.655 | 0.586 | 0.484 | 0.732 | 0.724 |
| **our-AllC**[†] | **0.662** | 0.805 | 0.770 | 0.514 | 0.672 | 0.626 | 0.655 | 0.736 | 0.752 |
| **our-NoM**[†] | 0.646 | 0.808 | 0.768 | 0.521 | 0.683 | 0.626 | 0.615 | 0.742 | 0.755 |
| **our**[†] | 0.658 | **0.812** | **0.774** | 0.522 | 0.685 | 0.633 | 0.604 | 0.741 | 0.757 |
| **model level transfer-learning**[†] | - | - | - | 0.531 | 0.684 | 0.637 | 0.645 | 0.784 | **0.795** |
| **feature level transfer-learning**[†] | - | - | - | 0.471 | 0.689 | 0.635 | **0.733** | **0.791** | 0.792 |

**Our-AllC**: 直接用未改进的loss
**Our-NoM**: 不用mention representation

**Model level transfer-learning**: 将模型先在Wiki训练数据上训练，把得到的参数用于在其他数据集上训练前的模型初始化

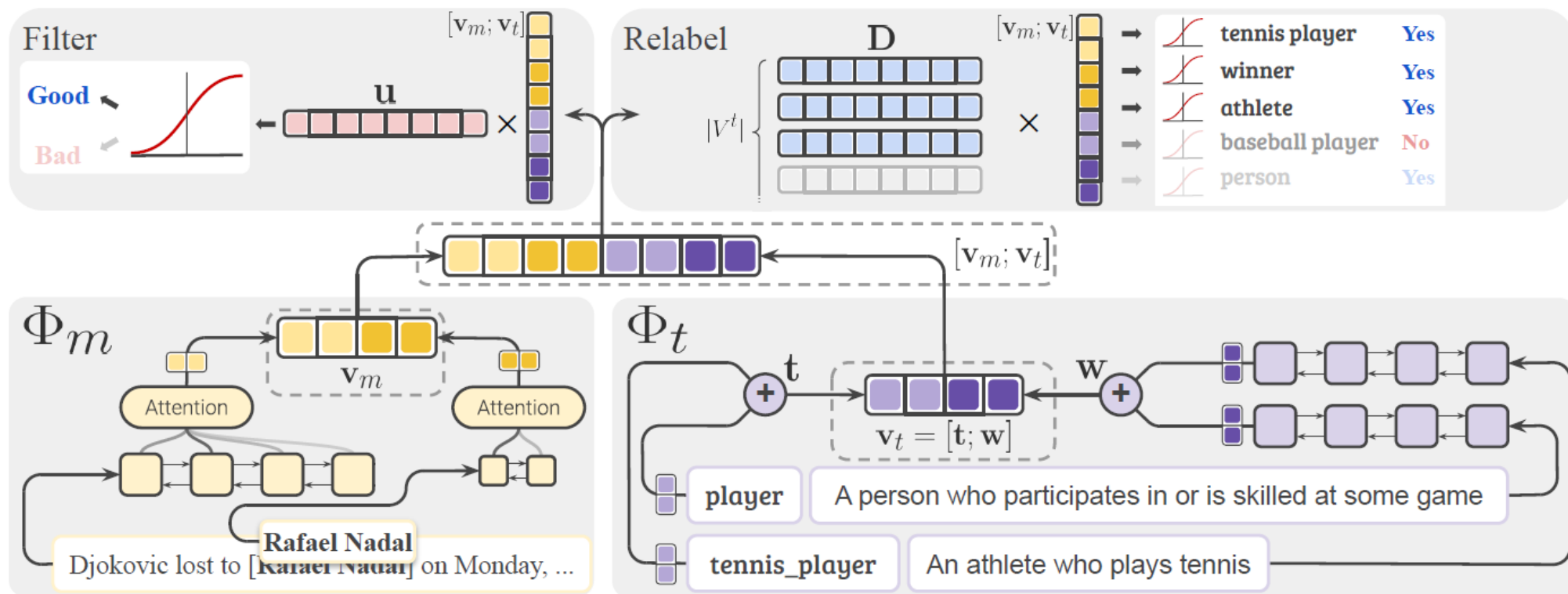**Feature level transfer-learning**: 将模型得到的feature representation用于另一个方法

# 应对弱监督训练数据

- 自动生成的训练数据标签不一定准确
- 许多缓解不准确标签不利影响的方法被提出，如
  - 改进训练loss
    - (Ren et al., 2016), (Abhishek et al., 2017)等
  - 对样本自动重标
    - (Onoe and Durrett, 2019)
  - 基于聚类的方法
    - (Chen et al., 2019)
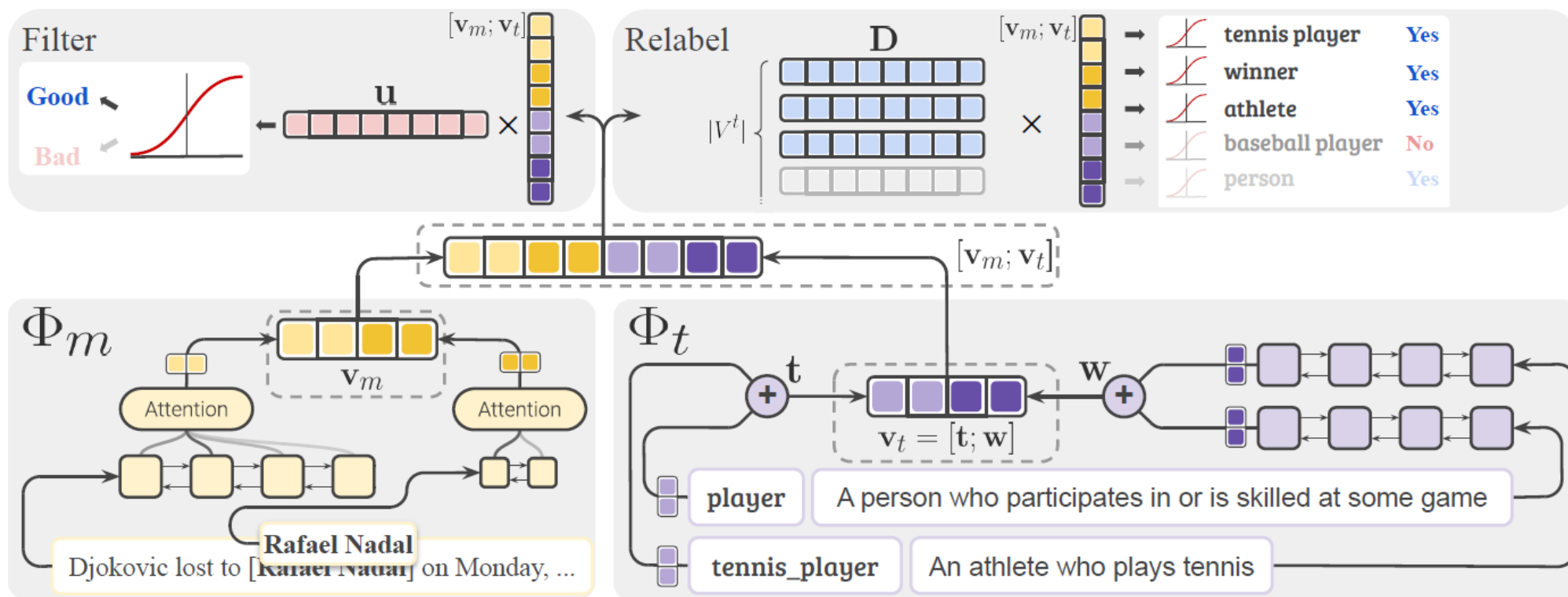  - 基于对抗学习的方法
    - (Shi et al., 2020)

# 应对弱监督训练数据

- 自动重标训练样本 (Onoe and Durrett, 2019)


- 设计一个可以自动重标训练样本的模型
  - 输入自动生成的样本（包括样本的标签也作为输入）
  - 模型目标是更正可能有错的标签，输出更准确的标签
- 用重标后的样本训练细粒度实体分类模型

Onoe, Y., & Durrett, G. (2019, June). Learning to Denoise Distantly-Labeled Data for Entity Typing. In *Proceedings of NAACL-HLT* (pp. 2407-2417).

# 自动重标训练样本



- 输入：一个弱监督样本，具体包含mention、上下文、弱监督标签、标签的定义
- Filter部分决定是否保留该样本（二分类）
- Relabel部分尝试对标签进行重标（多标签分类）

# 自动重标训练样本



**模型训练**：基于Ultrafine人工标注的2k个样本生成该模型的训练数据
如：对一个样本，人工标注了标签[person, athlete, player, tennis_player]，将标签改错为
[person, tennis_player, actor]就变成了一个弱监督样本

# 自动重标训练样本

- (Onoe and Durrett, 2019)的实验效果

| Model | P | R | F1 |
|---|---|---|---|
| Ours + GloVe w/o augmentation | 47.6 | 23.3 | 31.3 |
| Ours + ELMo w/o augmentation | **55.8** | 27.7 | 37.0 |
| Ours + ELMo w augmentation | 55.5 | 26.3 | 35.7 |
| Ours + ELMo w augmentation + filter & relabel | 51.5 | **33.0** | **40.2** |
| BERT-Base, Uncased | 51.6 | **33.0** | **40.2** |
| Choi et al. (2018) w augmentation | 47.1 | 24.2 | 32.0 |
| LABELGCN (Xiong et al., 2019) | 50.3 | 29.2 | 36.9 |

Ultrafine数据集上的效果（Macro-F1）

| Model | Acc. | Ma-F1 | Mi-F1 |
|---|---|---|---|
| Ours + ELMo w/o augmentation | 42.7 | 72.7 | 66.7 |
| Ours + ELMo w augmentation | 59.3 | 76.5 | 70.7 |
| Ours + ELMo w augmentation + filter & relabel | 63.9 | **84.5** | 78.9 |
| Ours + ELMo w augmentation by Choi et al. (2018) | **64.9** | **84.5** | **79.2** |
| BERT-Base, Uncased | 51.8 | 76.6 | 69.1 |
| Shimaoka et al. (2017) | 51.7 | 70.9 | 64.9 |
| AFET (Ren et al., 2016a) | 55.1 | 71.1 | 64.7 |
| PLE (Ren et al., 2016b) | 57.2 | 71.5 | 66.1 |
| Choi et al. (2018) | 59.5 | 76.8 | 71.8 |
| LABELGCN (Xiong et al., 2019) | 59.6 | 77.8 | 72.2 |

OntoNotes数据集上的效果

# 细粒度实体分类方法

• 基于预训练模型的方法

构建**BERT输入**：

[CLS] sentence [SEP] mention string [SEP]

输入**BERT**后得到[CLS]对应的向量表示，记为$u \in R^d$，基于$u$得到每个类别对应分数：

$$s = Wu$$

其中$W \in R^{K \times d}$，$K$为类别数，$s$中的每个元素对应一个类别

**损失函数**：margin-based loss或binary cross entropy loss

**实验结果**：

**Ultrafine**数据集：Precision 51.0 Recall 33.8 F1 40.7

**OntoNotes**数据集：Acc 63.25, Ma-F1 80.84, Mi-F1 75.90

Dai, Hongliang, Yangqiu Song, and Haixun Wang. "Ultra-Fine Entity Typing with Weak Supervision from a Masked Language Model." Proceedings of ACL. 2021.

# 基于大模型的方法

**HL** **You**

Given a sentence and an entity mention, output the fine-grained types of the entity mention in format ["type1", "type2", ...]

Sentence: FedEx is a major player in the package delivery market.
Entity Mention: FedEx

**ChatGPT**

In this context, "FedEx" is primarily an organization. The fine-grained types could include its industry or sector. Here's a potential output:

```json
["Organization", "Courier", "Delivery"]
```

# 其他细粒度实体分类方法

- 其他细粒度实体分类方法思路
  - 考虑类别体系中类别间的层次关系信息
  - 显示地利用上下文表示的语义关系
    - 如，文本中出现"$M$ and other $t$"，则意味着$t$ 可能是$M$的上位词
      - "... Microsoft and other companies ..."
  - 微调UFET模型为普通FET任务的模型
  - 等

Dai, Hongliang, and Ziqian Zeng. "From Ultra-Fine to Fine: Fine-tuning Ultra-Fine Entity Typing Models to Fine-grained." Proceedings of ACL. 2023.
Dai, Hongliang, Yangqiu Song, and Xin Li. "Exploiting semantic relations for fine-grained entity typing." AKBC. 2020.