



知识图谱

在线峰会

2021.03.27 (周六) 09:00~18:00





小米在知识表示学习 方向的探索与实践

吕荣荣 小米-算法工程师



目录

CONTENTS

01 业务介绍

小米知识图谱的架构和业务

03 算法应用

知识表示学习在实体链接、实体推荐、知识图谱补全的应用

02 算法介绍

融合文本和知识图谱的知识表示学习方法

04 总结

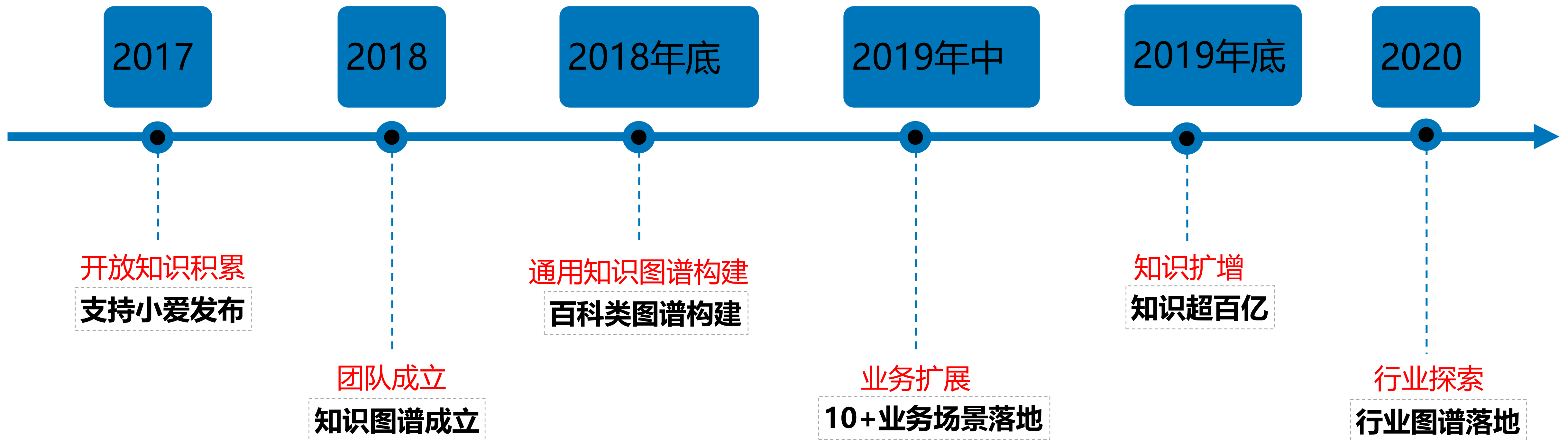
总结和展望

01

业务介绍

小米知识图谱的架构和业务

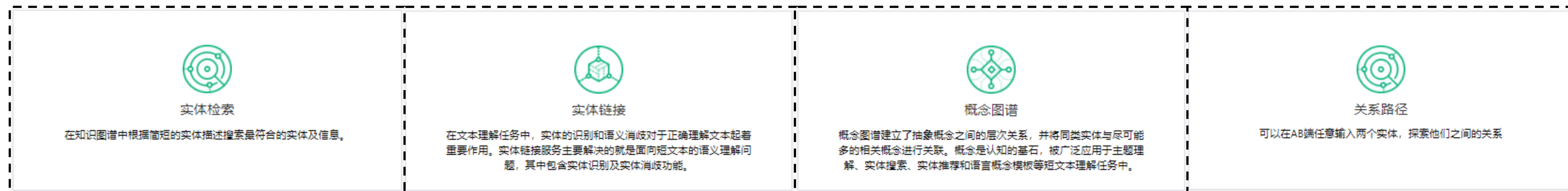
小米知识图谱的架构和业务



小米知识图谱的架构和业务



小米知识图谱团队，旨在研究知识图谱在开放领域和行业领域的构建和应用技术，把图谱推广到相关的业务场景上。



■ ■ ■ ■ ■

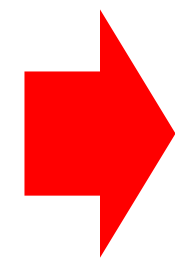
我们的用户



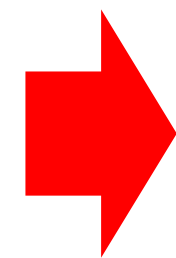
小米知识图谱的架构和业务



知识图谱团队为小爱同学赋能



Query:
巩俐的籍贯



下午5:42

人物·视频

巩俐的籍贯
籍贯：山东济南

《兰心大剧院》：巩俐情系中...

秒懂雷瓦特

图片

巩俐

简介:巩俐 (GongLi) , 1965年12月31日出生于中国辽宁省沈阳市, 祖籍山东省济南市, 华语电影女演员, 毕业于中央戏剧学院, 联合国促进和平艺术家, 联合国全球环境保护大使。1987年因主演电影《红高粱》成名, 该片获得第38届柏林国际电影节金熊奖。1992年凭借电影《秋菊打官司》获得第4...

籍贯: 山东济南

职业: 演员

出生日期: 1965年12月31日

个人信息: 169cm/56kg/摩羯座/A型血

本次回答满意吗? 展开更多

赞 踩

搜索

视频

图片

主要作品

霸王别姬

片名: 霸王别姬

首播: 1993-01-01(中国香港) /

评分: 9.6

主演: 张国荣、张丰毅、巩俐、葛优、英达

To Live

片名: 活着

首播: 1994-05-17(戛纳电影节) /

评分: 9.2

主演: 葛优、巩俐、姜武、牛犇、郭涛

唐伯虎点秋香

片名: 唐伯虎点秋香

首播: 1993-07-01(香港)

评分: 8.5

主演: 周星驰、巩俐、陈百祥、郑佩佩、朱咪咪

艺伎回忆录

片名: 艺伎回忆录

首播: 2005-12-10(日本)

人物关系

让·米歇尔·雅尔

丈夫

赵英

母亲

巩俐

哥哥

黄和祥

前夫

周星驰

搭档

刘德华

搭档

张国荣

搭档

张艾嘉

搭档

张艺谋

搭档

相关资讯

进入

巩俐,别再整那些烂活儿了

要说华人女星,绕不开一个人,巩俐。江湖地位不用多说,“巩皇”这名头不是白叫的。但近年我们...

离开张艺谋之后,为何章子怡强过巩俐?“天...

说到最强“谋女郎”的比拼,一定是在巩俐和章子怡之间产生,这个因为和张艺谋合作而诞生的女演员称号,是张艺谋整个职业生涯的重要标志。当年...

巩俐领衔顶配阵容《夺冠》提档9月25日

影片由陈可辛执导,张冀编剧,巩俐、黄渤、吴刚、彭昱畅、白浪、中国女子排球队领衔主演,李现特别出演。影片《夺冠》原定1月24日上映,1月19...

为什么张艺谋不娶巩俐,而娶了小他31岁的...

阅读全文约需8分钟为什么张艺谋不娶巩俐,而娶了小他31岁的陈婷?文/晏凌羊 来源/晏凌羊每次看到...

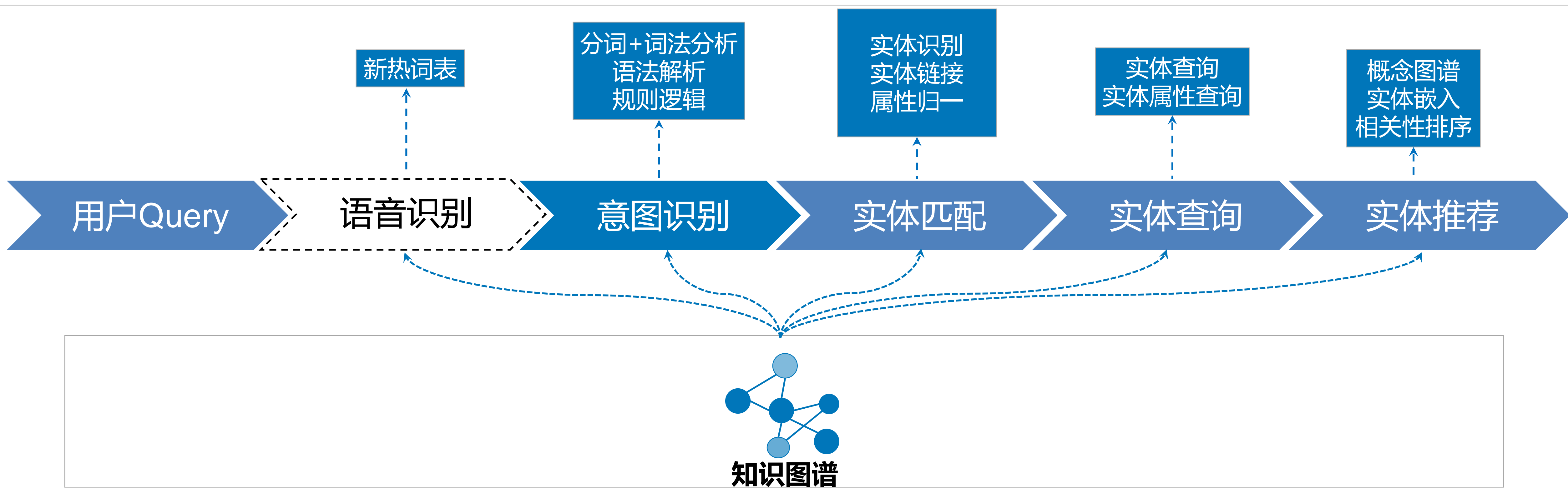
上映27天拿下23天日冠,八佰之后,巩俐这部...

而《夺冠》这部由巩俐主演的电影,虽然经历了改

小米知识图谱的架构和业务



知识图谱团队为小爱同学赋能



For Example:

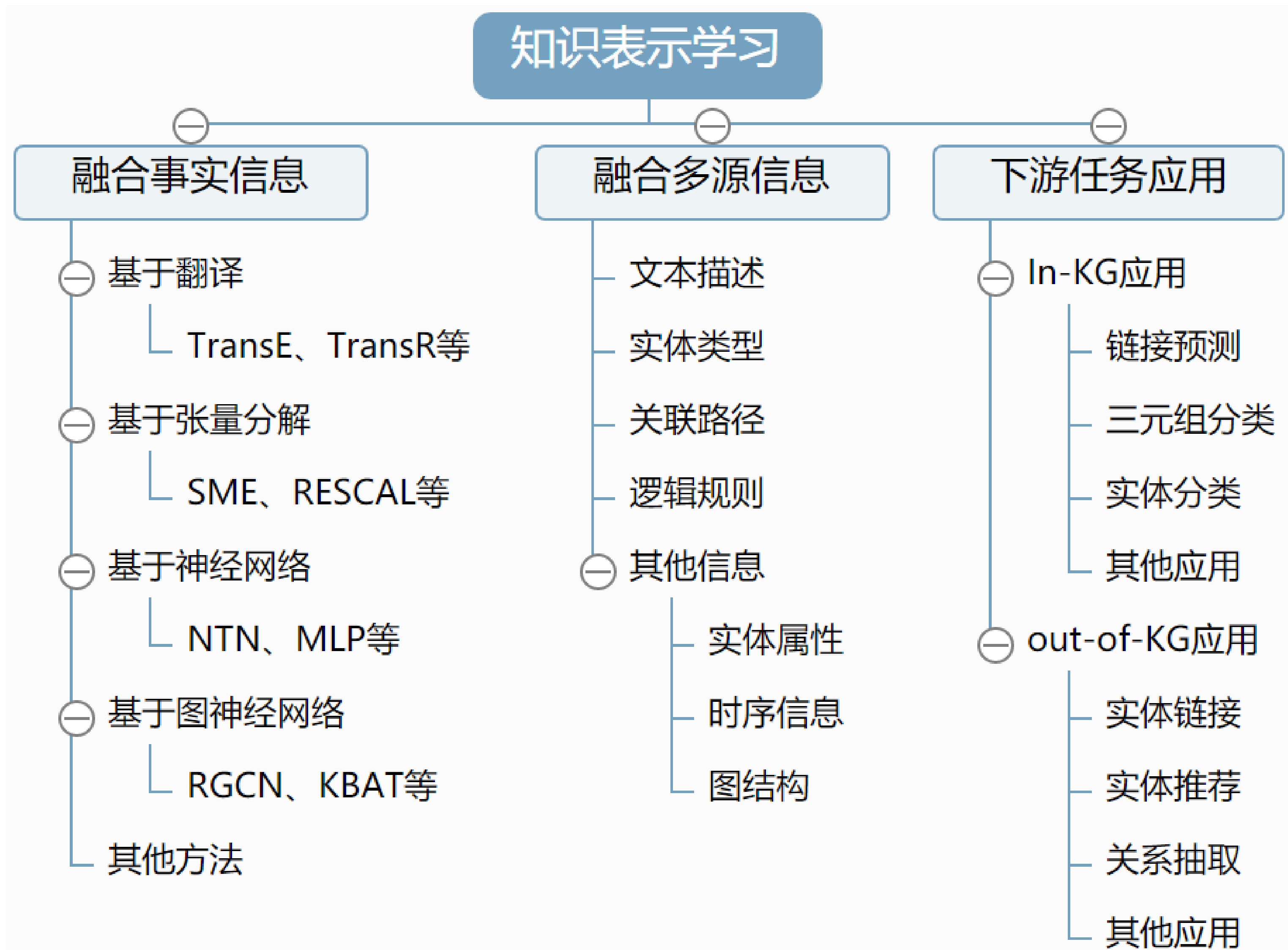


02

算法介绍

融合多源信息 的知识表示学习 方法

融合文本和知识图谱的知识表示学习方法



问题：

仅利用知识图谱的三元结构信息进行表示学习，数据稀疏问题严重，尚有大量与知识有关的其他信息没有得到有效利用。

解决方法：

1. 利用知识库中的其他信息，如实体和关系的描述信息等。
2. 利用知识库外的海量信息，如互联网文本等包含大量与知识库实体和关系有关的信息。

融合文本和知识图谱的知识表示学习方法



融入文本描述的优势:

1. 可以发掘实体间的语义相关性, 精确的语义表述能够提升三元组的可区分性。
2. 可以解决zero-shot问题

实体的描述文本:

将所有三元组的“属性-属性值”或“关系-实体提及”都拼成一个字符串, 当作该实体的文本描述。由于 type 字段, 义项描述和摘要字段的信息更重要, 描述文本中都按照 type、义项描述、摘要和其他三元组的顺序进行拼接。

类型: 人物|描述: 中国中央电视台节目主持人|简介: 撒贝宁, 1976年3月23日出生于广东省湛江市, 籍贯湖北武汉, 祖籍安徽和县.....

类型: 人物|描述: 撒贝宁妻子|简介: 李白 (Lisa), 加拿大人, 北京大学教育学博士, 曾经是“五洲唱响”组合成员, 当年该团体在星.....



类型: 人物|描述: 手游《王者荣耀》中的英雄角色|简介: 李白, 是腾讯手游《王者荣耀》中的一位刺客型英雄角色, 原型为唐代诗人“诗仙”.....

类型: 游戏|描述: 2015年腾讯天美发行的MOBA手游|简介: 《王者荣耀》是由腾讯游戏天美工作室群开发并运行的一款运营在.....

融合文本和知识图谱的知识表示学习方法

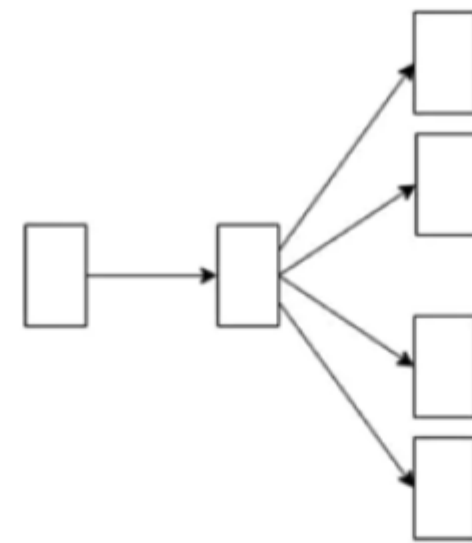


Text Model: Skip-Gram

$$z(w, v) = b - 0.5 \cdot \|\mathbf{w} - \mathbf{v}\|_2^2$$

$$\Pr(w|v) = \frac{\exp\{z(w, v)\}}{\sum_{\tilde{w} \in \mathcal{V}} \exp\{z(\tilde{w}, v)\}}$$

$$\mathcal{L}_T = - \sum_{(w, v)} \log \Pr(w|v)$$

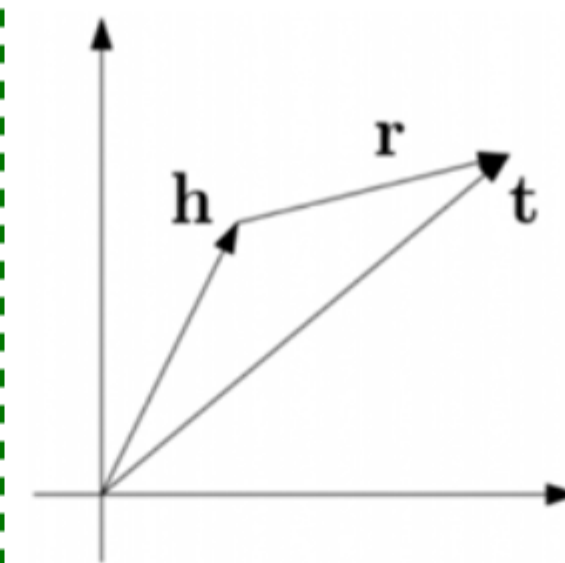


Knowledge Model: TransE

$$z(h, r, t) = b - 0.5 \cdot \|\mathbf{h} + \mathbf{r} - \mathbf{t}\|_2^2$$

$$\Pr(h|r, t) = \frac{\exp\{z(h, r, t)\}}{\sum_{\tilde{h} \in \mathcal{I}} \exp\{z(\tilde{h}, r, t)\}}$$

$$\mathcal{L}_K = - \sum_{(h, r, t)} [\log \Pr(h|r, t) + \log \Pr(t|h, r) + \log \Pr(r|h, t)]$$



实体描述

关系三元组

联合表示学习

知识向量表示

$$z(e, w) = b - 0.5 \cdot \|\mathbf{e} - \mathbf{w}\|_2^2$$

$$\Pr(w|e) = \frac{\exp\{z(e, w)\}}{\sum_{\tilde{w} \in \mathcal{V}} \exp\{z(e, \tilde{w})\}}$$

$$\Pr(e|w) = \frac{\exp\{z(e, w)\}}{\sum_{\tilde{e} \in \mathcal{E}} \exp\{z(\tilde{e}, w)\}}$$

$$\mathcal{L}_A = - \sum_{e \in \mathcal{E}} \sum_{w \in D_e} [\log \Pr(w|e) + \log \Pr(e|w)]$$

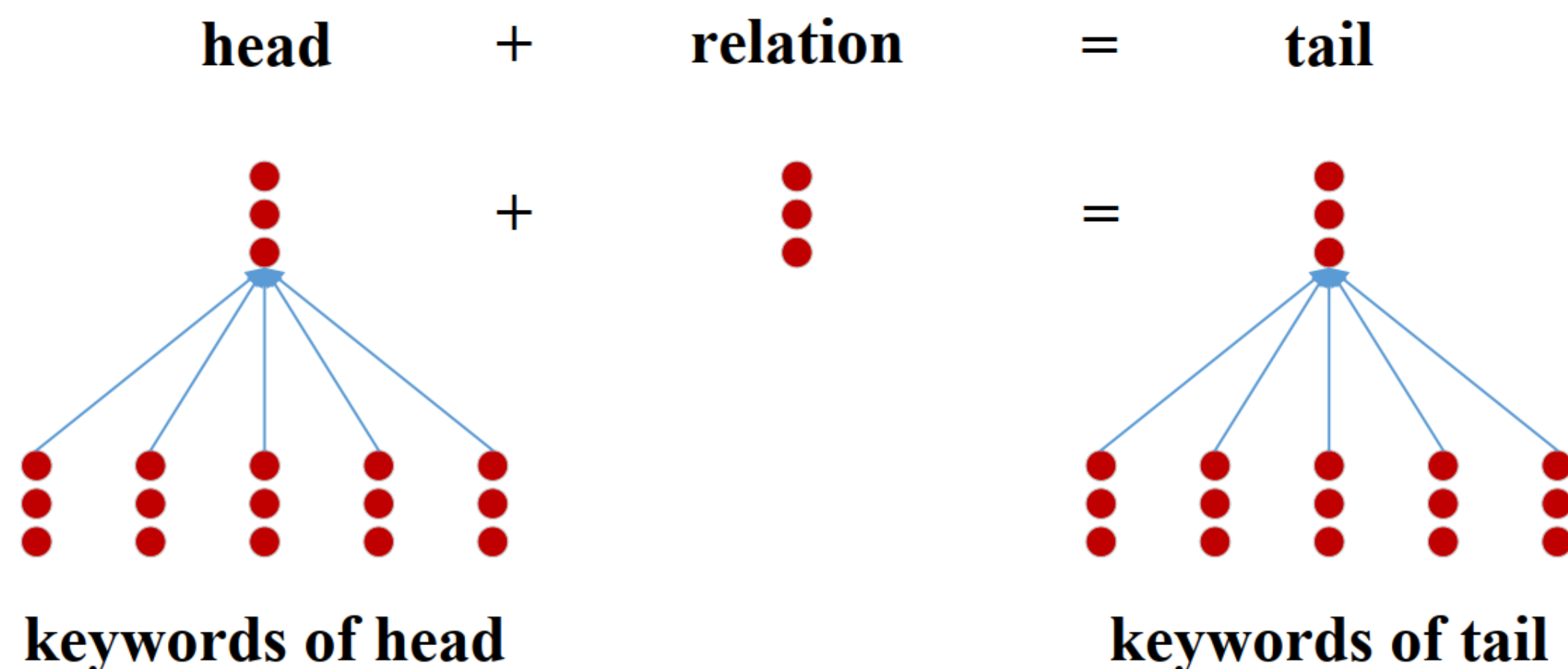
$$\mathcal{L}(\{\mathbf{e}_i\}, \{\mathbf{r}_j\}, \{\mathbf{w}_l\}) = \mathcal{L}_K + \mathcal{L}_T + \mathcal{L}_A$$

1. 第一部分是文本嵌入模型，采用Skip-Gram模型，距离公式采用的是两个单词之间的欧氏距离。
2. 第二部分是知识嵌入模型，采用transE模型。
3. 第三部分是对齐模型，利用文本描述对齐，确保实体，关系能够和文本中的单词在同一个语义空间中。

融合文本和知识图谱的知识表示学习方法



DKRL(CBOW): The CBOW Encoder



$$E_S = || h_s + r + h_s ||$$

$$E_{DD} = || h_d + r - t_d ||$$

$$E_{DS} = || h_d + r - t_s ||$$

$$E_{SD} = || h_s + r - t_d ||$$

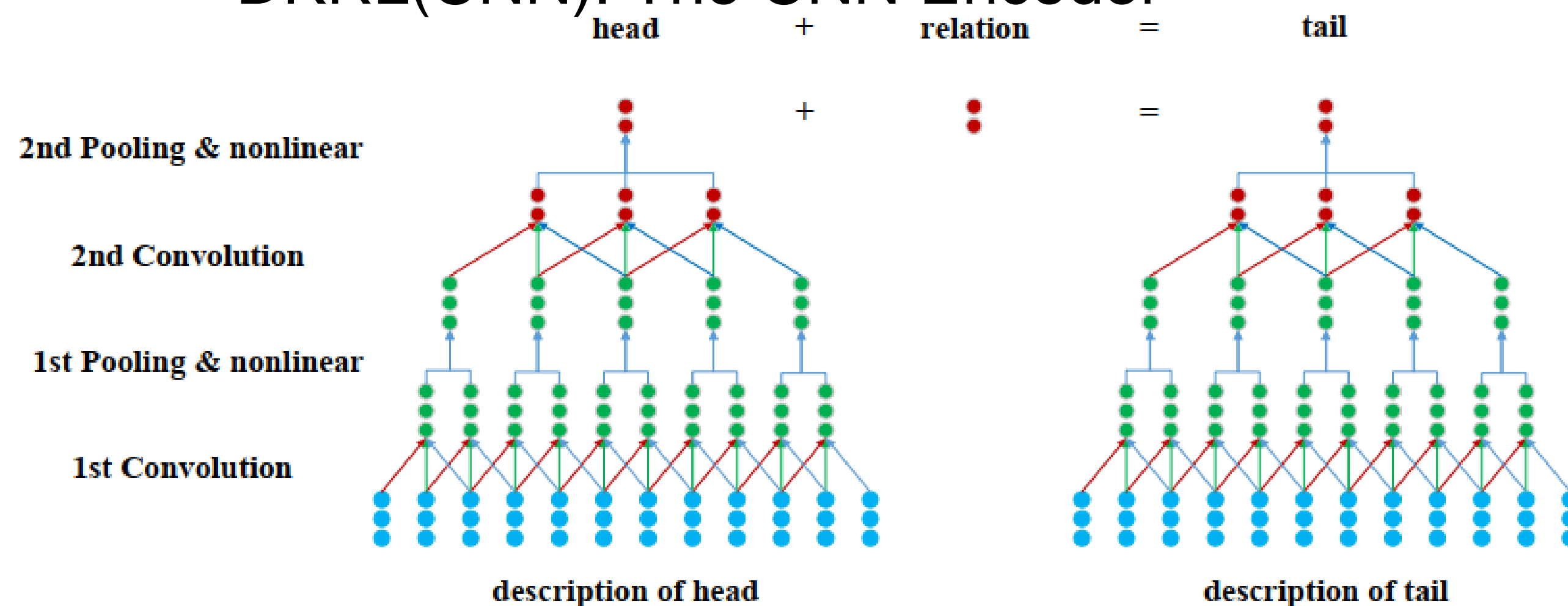
$$E = E_S + E_D$$

$$= E_S + E_{DD} + E_{DS} + E_{SD}$$

$$L = \sum_{(h,r,t) \in T} \sum_{(h',r',t') \in T'} \max(\gamma + d(h + r, t) - d(h' + r', t'), 0),$$

$$T' = \{(h', r, t) | h' \in E\} \cup \{(h, r, t') | t' \in E\} \cup \{(h, r', t) | r' \in R\}$$

DKRL(CNN): The CNN Encoder



融合文本和知识图谱的知识表示学习方法



需求目标:

- 1. 得到实体向量
- 2. 得到词向量
- 3. 词向量和实体向量之间可以计算相似度
- 4. 实体向量和实体向量之前可以计算相似度

Table 1: Mean Rank and HITS@10 of Knowledge Graph Completion (For Predicting Entity) on FB15K

FB15K	Mean Rank		HITS@10	
TransE	210	119	48.5	66.1
TransH	212	87	45.7	64.4
Jointly	167 ¹	39 ¹	51.7 ¹	77.3 ¹
DKRL(BOW)	200	113	44.3	57.6
DKRL(ALL)	181	91	49.6	67.4

词向量与实体向量的相似度

```
input word: 王者荣耀
entity and score:
李白 手游《王者荣耀》中的英雄角色 0.6871914871010902
李白 漫画《尸兄》中的角色 0.4868442647918333
李白 QQ游戏《英雄杀》卡牌 0.47989461648584353
李白 南宋徐钧所写诗歌 0.46257176254889987
李白 李荣浩演唱歌曲 0.4526289976719955
李白 撒贝宁妻子 0.4438043542986664
李白 唐代诗人 0.35362615078817616
李白 北京人艺经典话剧 0.34886872301725935
李白 中共党员，上海地下党联络员 0.2934343603235514
```

```
input word: 将进酒
entity and score:
李白 唐代诗人 0.7954981104666751
李白 南宋徐钧所写诗歌 0.6594177266802413
李白 李荣浩演唱歌曲 0.5860849811810142
李白 手游《王者荣耀》中的英雄角色 0.5421541413927163
李白 北京人艺经典话剧 0.5303117604279298
李白 QQ游戏《英雄杀》卡牌 0.4655070254651763
李白 漫画《尸兄》中的角色 0.43050382437348667
李白 撒贝宁妻子 0.4236790438669058
李白 中共党员，上海地下党联络员 0.2748187908448527
```

实体向量与实体向量的相似度

```
input entity: 小米 禾本科狗尾草属一年生草本
entity and score:
苹果 蔷薇科苹果属植物 0.8354706297404012
苹果 伊朗1998年莎米拉·玛克玛尔巴夫执导电影 0.63257975713
苹果 韩国2008年康理贯执导电影 0.5702162436015609
苹果 安与骑兵演唱歌曲 0.5607941649198853
苹果公司 0.5381222371491942
苹果 网游《天堂梦》中人物 0.5211885607135783
苹果 2007年李玉执导电影 0.47746202503045454
苹果 动漫《男子高中生的日常》中角色 0.47035311617841247
```

```
input entity: 北京小米科技有限责任公司
entity and score:
苹果公司 0.79185107258104
苹果 安与骑兵演唱歌曲 0.6309651542813589
苹果 伊朗1998年莎米拉·玛克玛尔巴夫执导电影 0.5943178076
苹果 韩国2008年康理贯执导电影 0.5558840673715467
苹果 网游《天堂梦》中人物 0.533240821748632
苹果 蔷薇科苹果属植物 0.4861939169898948
苹果 2007年李玉执导电影 0.45654442486800906
苹果 动漫《男子高中生的日常》中角色 0.3878173304782322
```


03

算法应用



知识表示学习在小米的应用

- a) 实体链接
- b) 实体推荐
- c) 知识补全

实体链接中的应用



实体链接，就是把文本中的实体指称链接到知识图谱对应的实体上的任务。

Knowledge Graph (知识图谱)：一种语义网络，旨在描述客观世界的概念实体及其之间的关系。

Mention (提及、实体指称)：自然文本中表达实体的语言片段。

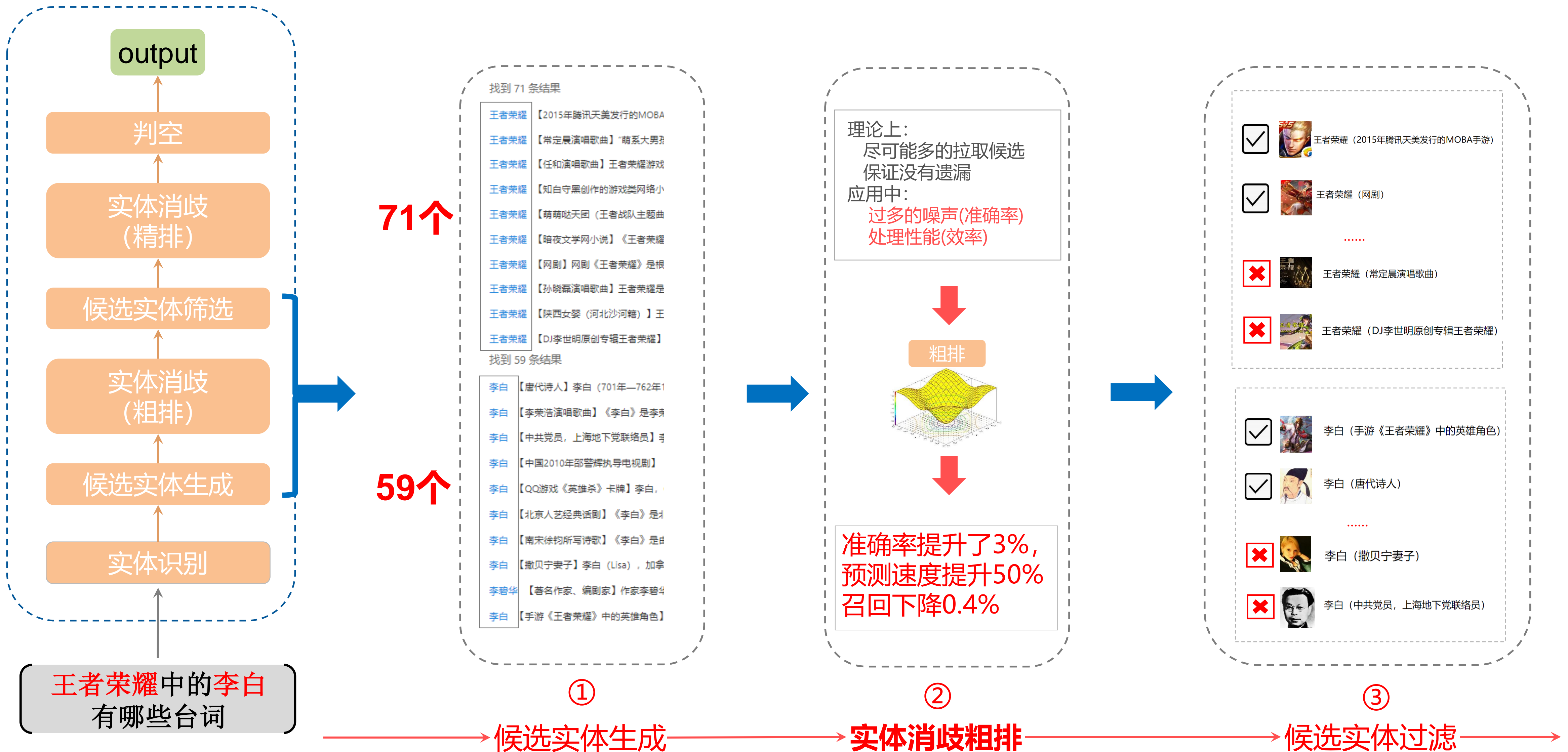
Entity (实体)：实体是知识库的基本单元，也是文本中承载信息的重要语言单位，由多个三元组组成



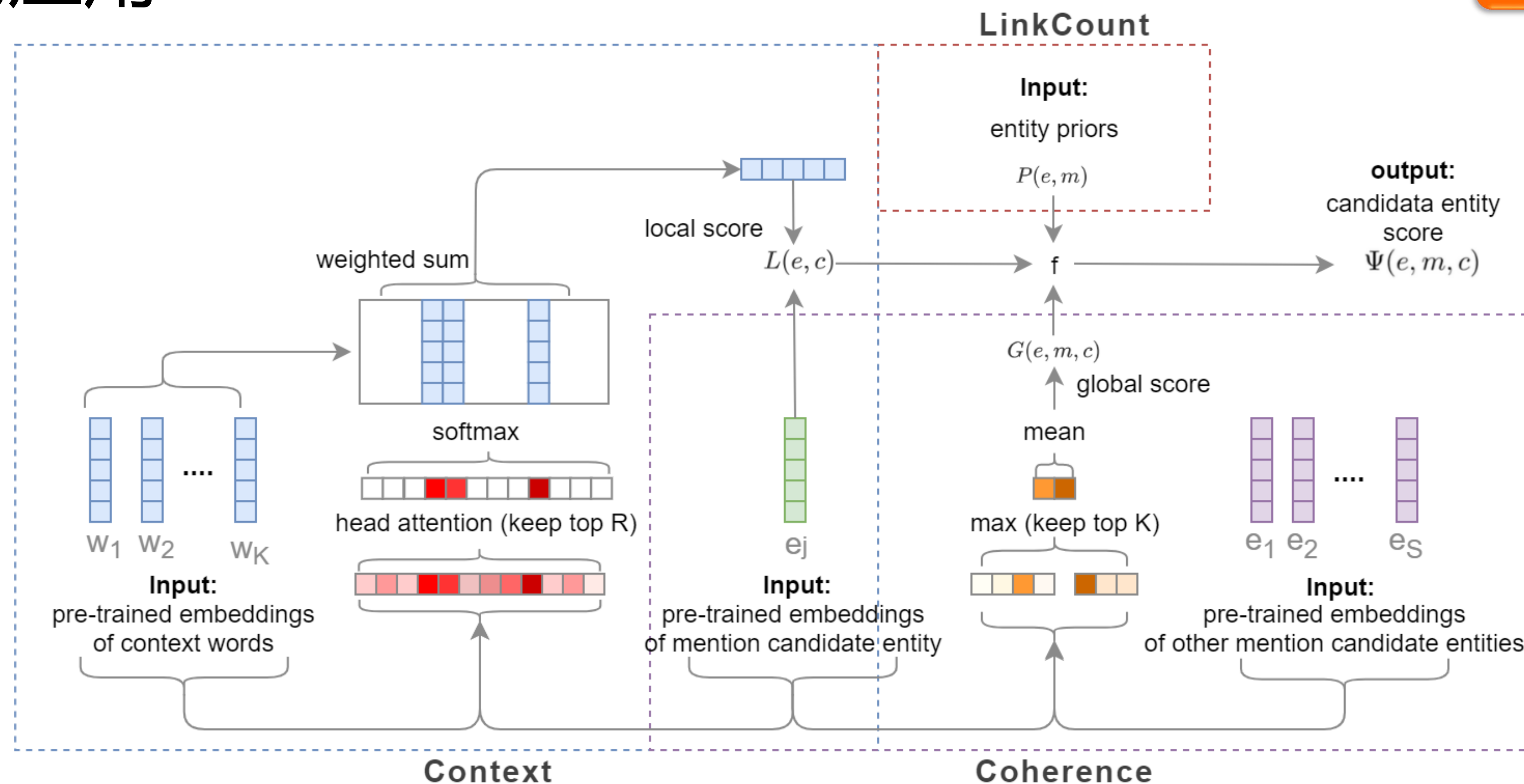
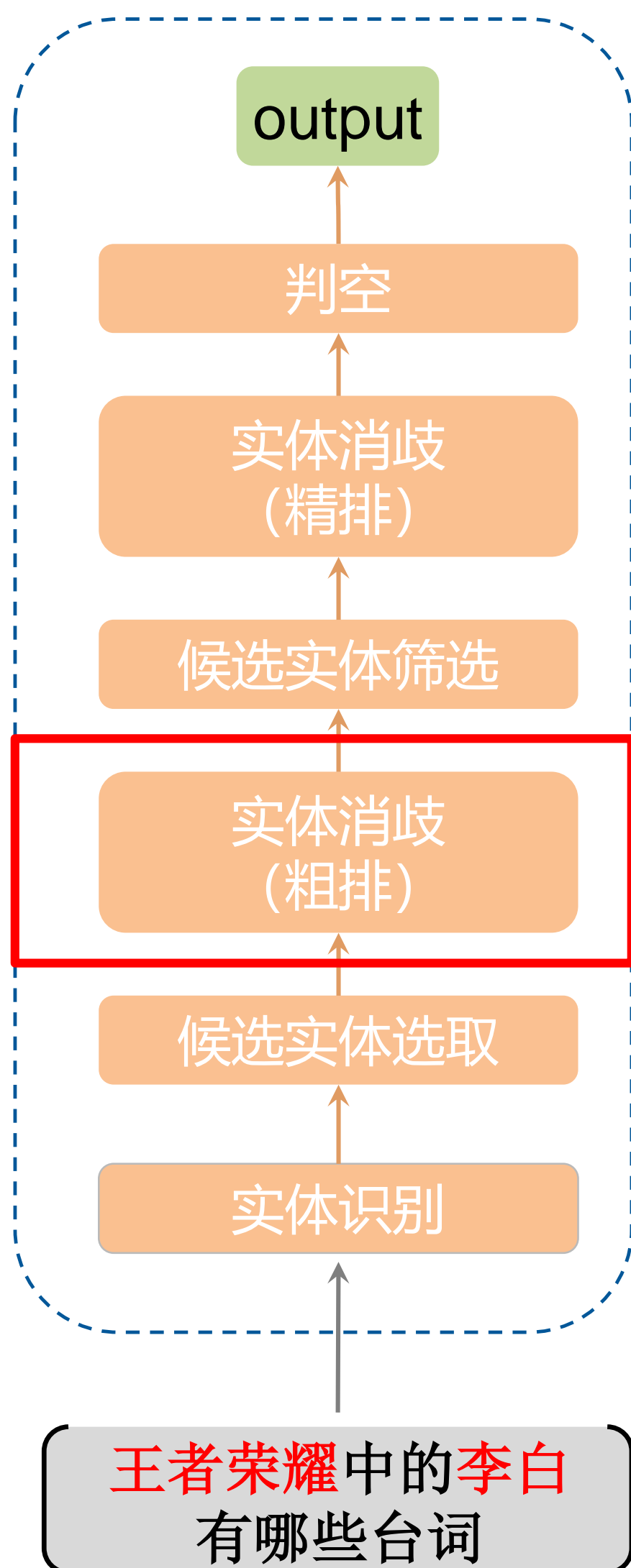
难点: Mention Variations: 同一entity有不同的mention。 (<李白>: 青莲居士、李太白)

Entity Ambiguity: 同一mention对应不同的entity。 ("李白": 王者荣耀中李白的技能; 李白和杜甫并称为为什么?)

实体链接中的应用



实体链接中的应用



$$u(w) = e^T X_w$$

$$\bar{c} = \{w \in c | u(w) \in \text{topR}(\mathbf{u})\}$$

$$\beta(w) = \begin{cases} \frac{\exp[u(w)]}{\sum_{v \in \bar{c}} \exp[u(v)]} & \text{if } w \in \bar{c} \\ 0 & \text{otherwise} \end{cases}$$

$$L(e, c) = \sum_{w \in \bar{c}} \beta(w) e^T X_w$$

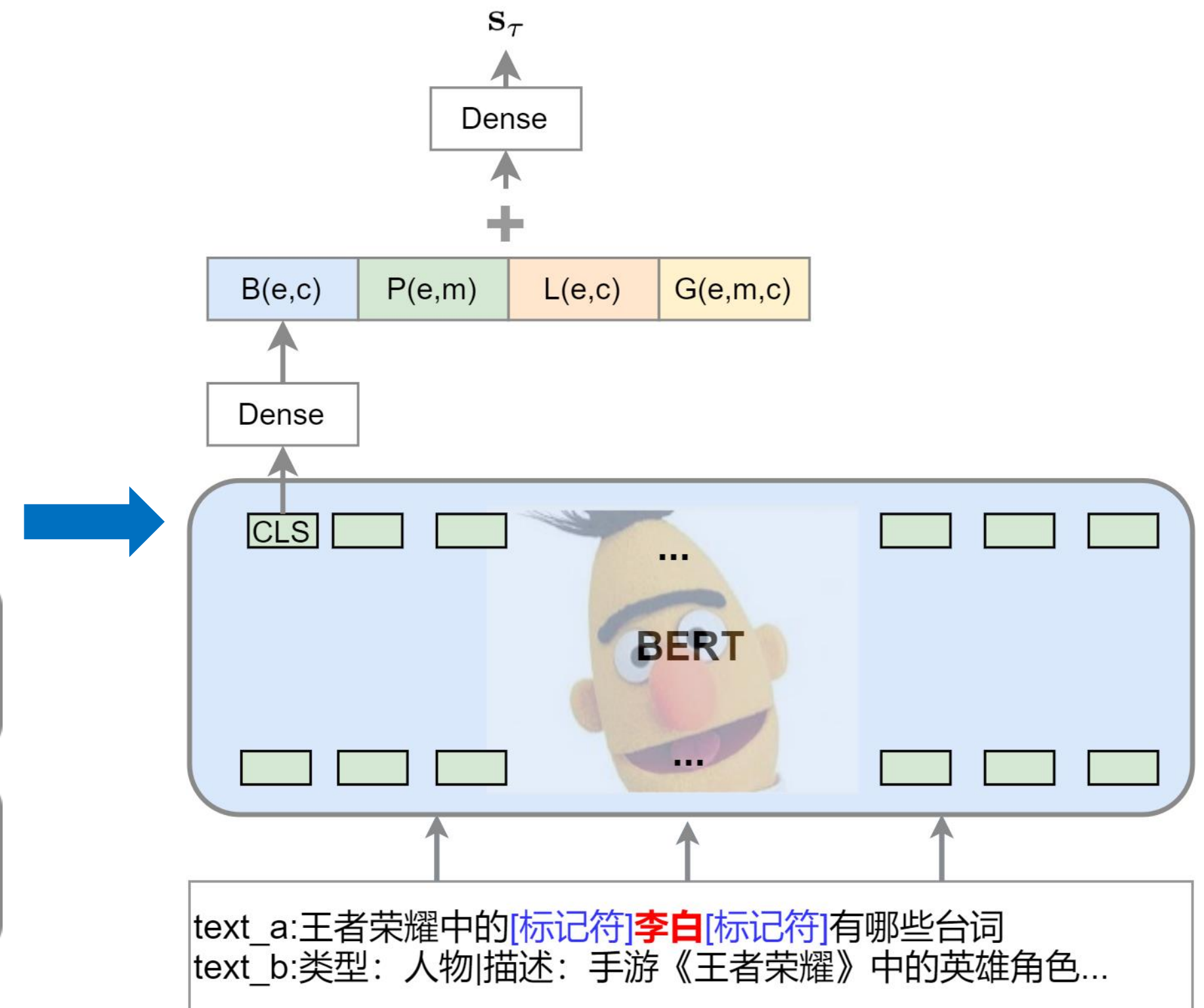
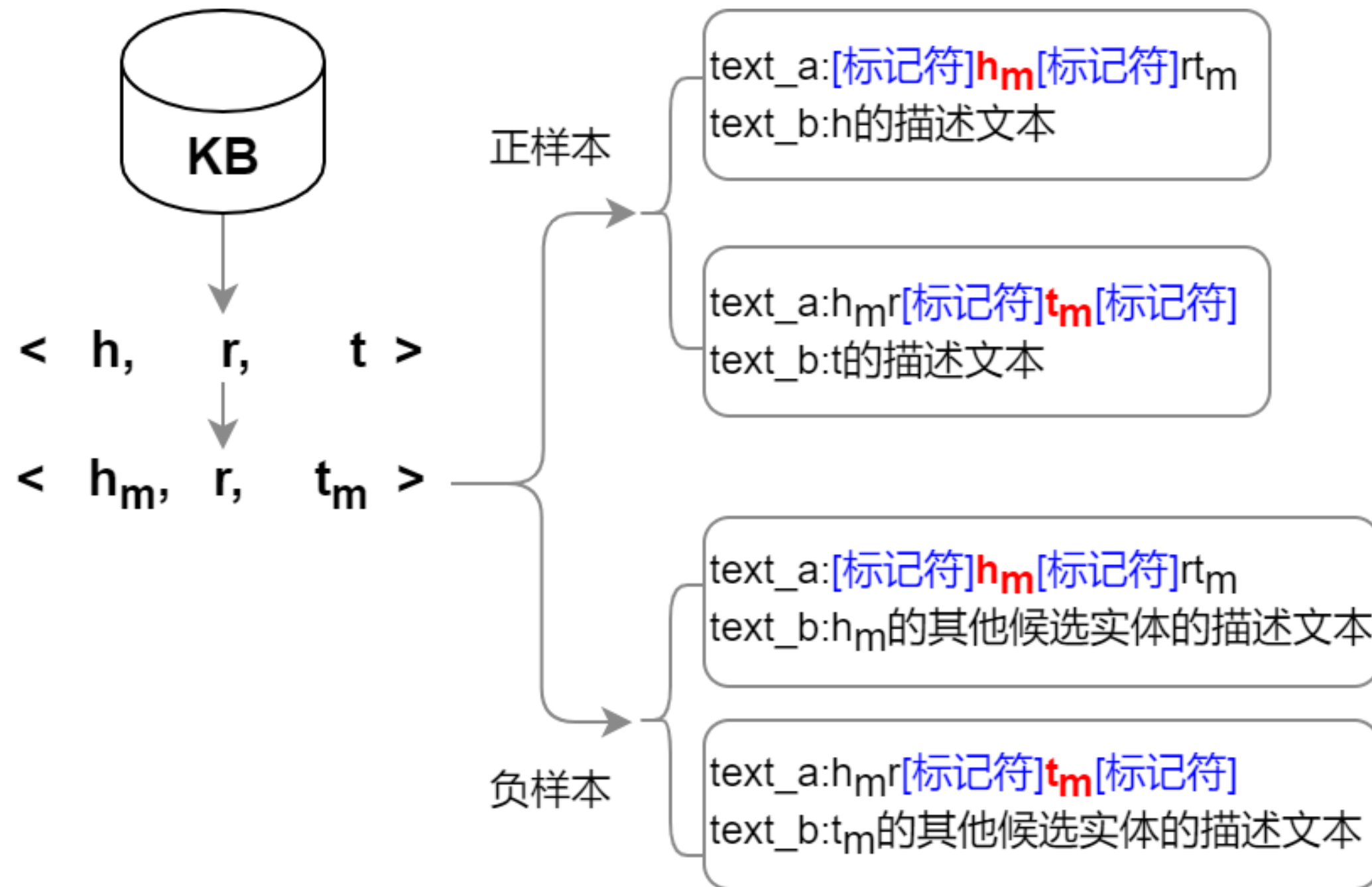
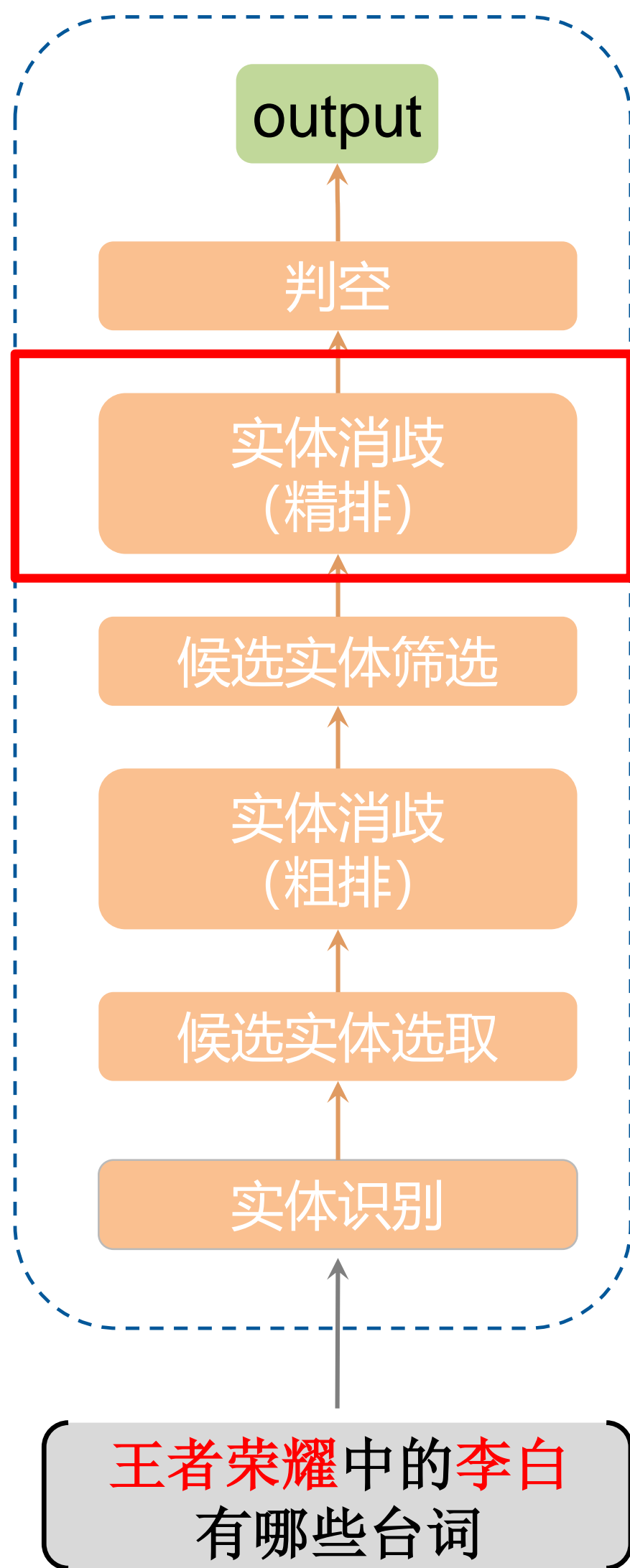
$$M_K = \{\bar{m} \mid \bar{m} \in c \wedge \bar{m} \neq m\}$$

$$G(e, m, c) = \frac{\sum_{\bar{m} \in M_K} \max(e^T X_{\bar{m}})}{K}$$

$$P(e, m) = \frac{\text{count}_m(e)}{\sum_{e \in X_e} \text{count}_m(e)}$$

$$\Psi(e, m, c) = f(P(e, m), L(e, c), G(e, m, c))$$

实体链接中的应用



二次训练:

第一次基于KB中的关系三元组训练
第二次基于标注的线上query训练

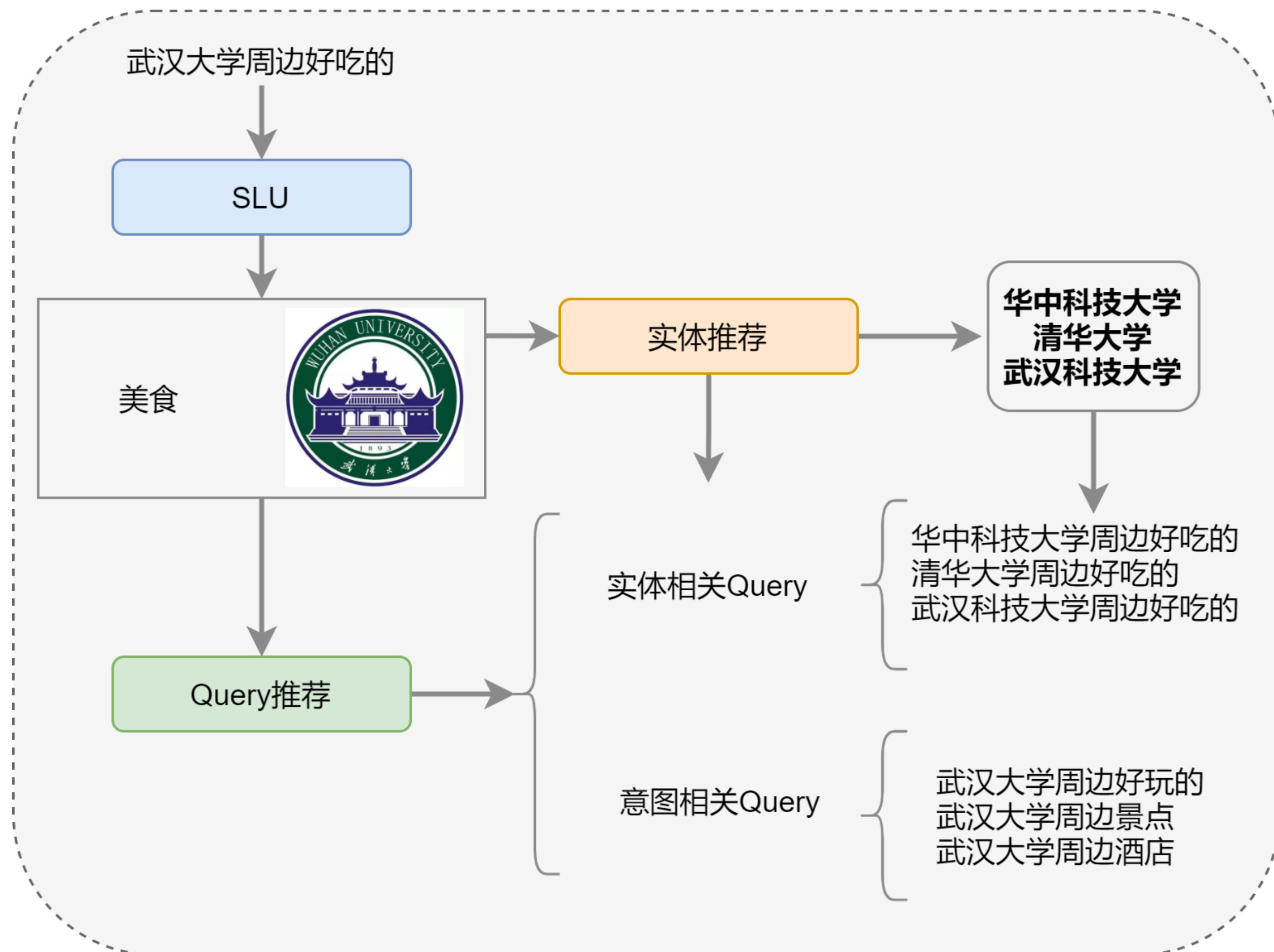
实体的描述文本: 将所有三元组的“属性-属性值”或“关系-实体提及”都拼成一个字符串, 当作该实体的文本描述。由于 type 字段, 义项描述和摘要字段的信息更重要, 描述文本中都按照 type、义项描述、摘要和其他三元组的顺序进行拼接。

$$s_{\tau} = f(h, r, t) = \text{sigmoid}(CW^T)$$
$$\mathcal{L} = - \sum_{\tau \in \mathbb{D}^+ \cup \mathbb{D}^-} (y_{\tau} \log(s_{\tau 0}) + (1 - y_{\tau}) \log(s_{\tau 1}))$$

实体推荐中的应用



实体推荐，就是根据给定的实体，推荐一系列相关实体。



实体推荐中的应用



实体推荐冷启动问题：如何寻找实体的相关实体？

方案：通过百科页面、关系三元组、经过实体链接处理的新闻中的共现实体对，经过类别过滤，作为正样本进行训练。

小米集团

中文维基百科【维基百科中文网站】
(重定向自小米科技)

此条目介绍的是中国大陆的科技企业。关于一种粮食，请见“小米”。

小米集团（英语：Xiaomi Corporation，简称：**小米**，**港交所**：**1810**）是中国一家从事智能硬件和电子产品研发、智能家居生态建设的大型移动互联网企业，成立于2010年4月6日^[9]，总部位于中国北京。截至2019年，小米集团在职员工人数近1.82万^[10]，并在全球超过90个国家和地区的市场开展业务。2018年7月9日，小米以“小米集团”名义在**香港交易所**主板挂牌上市，成为港交所上市制度改革后首家采用不同投票权架构的上市企业。^[11]小米还是继**苹果**、**三星**、**华为**之后第四家拥有手机芯片自研能力的手机公司。

通过旗下生态链品牌MIJIA（米家）与旗下子品牌Redmi（红米）、POCO，小米的产品线从智能手机及耳机、移动电源等手机周边产品和**音箱**、手环等相关移动智能硬件，扩展到智能电视、**机顶盒**、**路由器**、**空气净化器**、**电饭煲**等家居消费产品。^[12]

小米已建成全球最大消费类IoT物联网平台，连接超过2.52亿台智能设备（不含智能手机和个人电脑）^[13]，MIUI月活跃用户达到3.096亿^[14]。小米系投资的公司接近400家，覆盖智能硬件、生活消费用品、教育、游戏、社交网络、文化娱乐、医疗健康、汽车交通、金融等领域。

自2019年，小米连续入选**世界500强**排行榜和“BrandZ全球最具价值品牌百强”^[15]。2019年，小米首次入选世界500强，排名468位^[16]，2020年排名第422位。小米于2019年8月在中国互联网协会、工信部信息中心发布的2019年中国互联网企业100强中排名15^[17]。2019年10月，在**福布斯**发布的《2019**福布斯**全球数字经济100强榜》位列第56位。在2019《**财富**》未来50强榜单中排名第7^[18]。

香港财经界把**阿里巴巴**（**港交所**：**9988**）、**腾讯**（**港交所**：**700**）、**美团点评**（**港交所**：**3690**）、小米四只中国大陆科技股的英文名称首字母，合称“ATMX”股份。^[19]



相关资讯

进入

巩俐,别再整那些烂活儿了

要说华人女星,绕不开一个人,巩俐。江湖地位不用多说,“巩皇”这名头不是白叫的。但近年我们...

离开张艺谋之后,为何章子怡强过巩俐?“天...

说到最强“谋女郎”的比拼,一定是在巩俐和章子怡之间产生,这个因为和张艺谋合作而诞生的女演员称号,是张艺谋整个职业生涯的重要标志。当年...

巩俐领衔顶配阵容《夺冠》提档9月25日

影片由陈可辛执导,张冀编剧,巩俐、黄渤、吴刚、彭昱畅、白浪、中国女子排球队领衔主演,李现特别出演。影片《夺冠》原定1月24日上映,1月19...

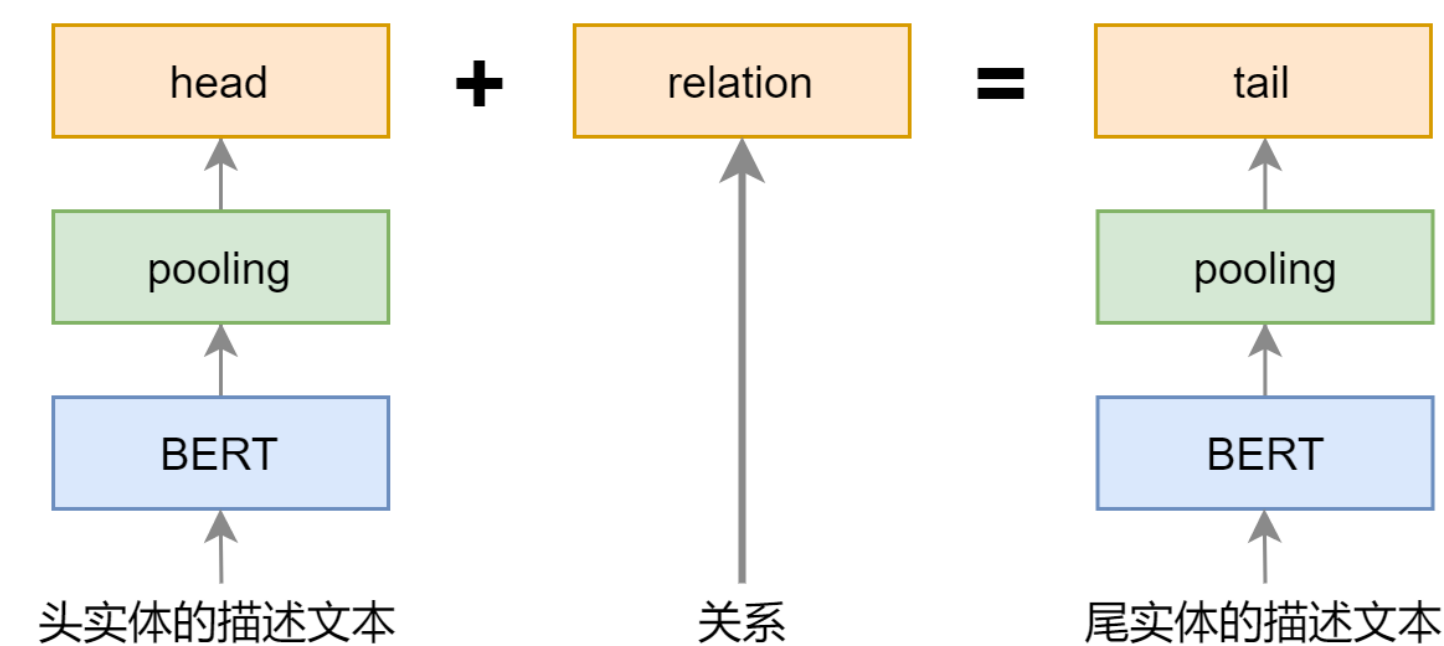
为什么张艺谋不要巩俐,而娶了小他31岁的...

阅读全文约需8分钟为什么张艺谋不要巩俐,而娶了小他31岁的陈婷?文/晏凌羊 来源/晏凌羊每次看到...

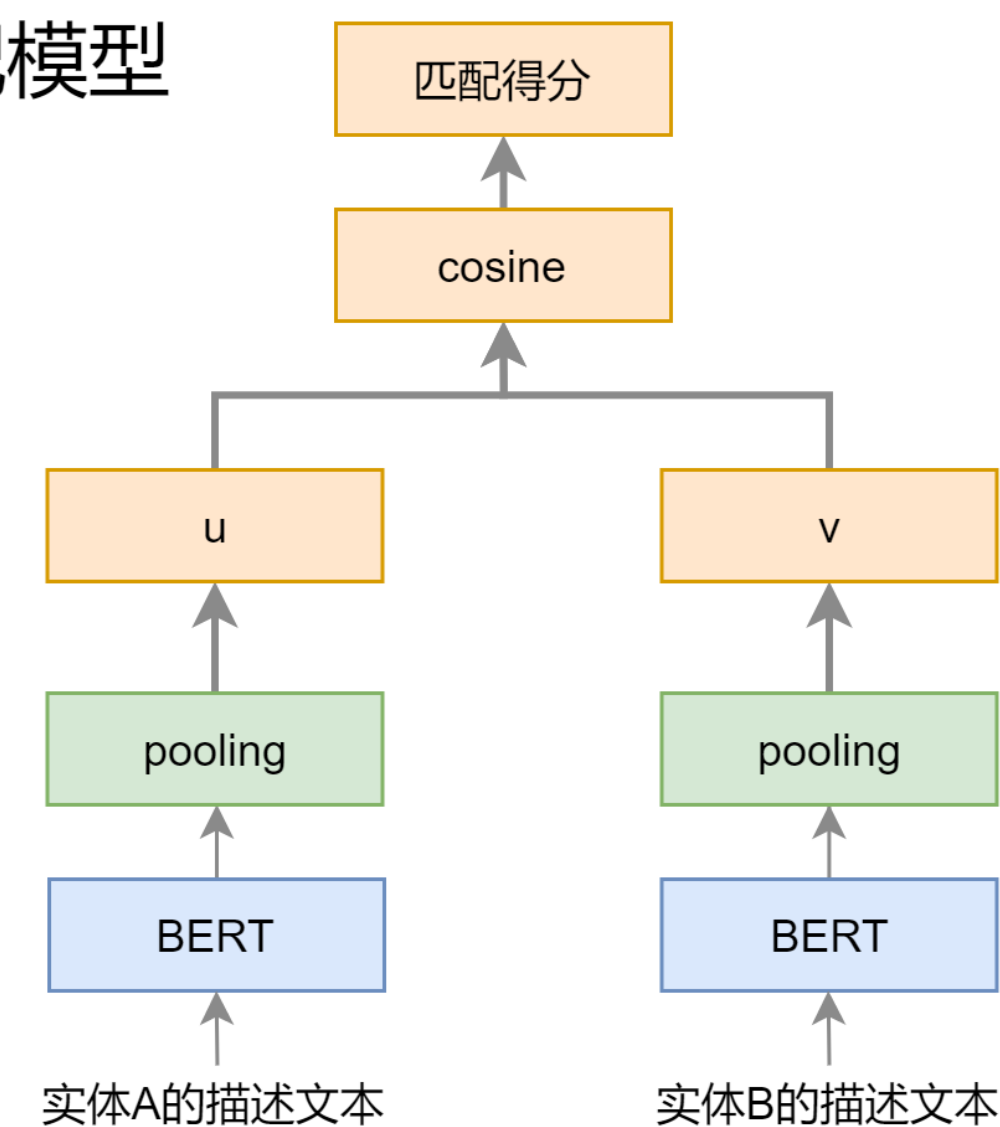
上映27天拿下23天日冠,八佰之后,巩俐这部...

而《夺冠》这部由巩俐主演的电影,虽然经历了改

表示模型



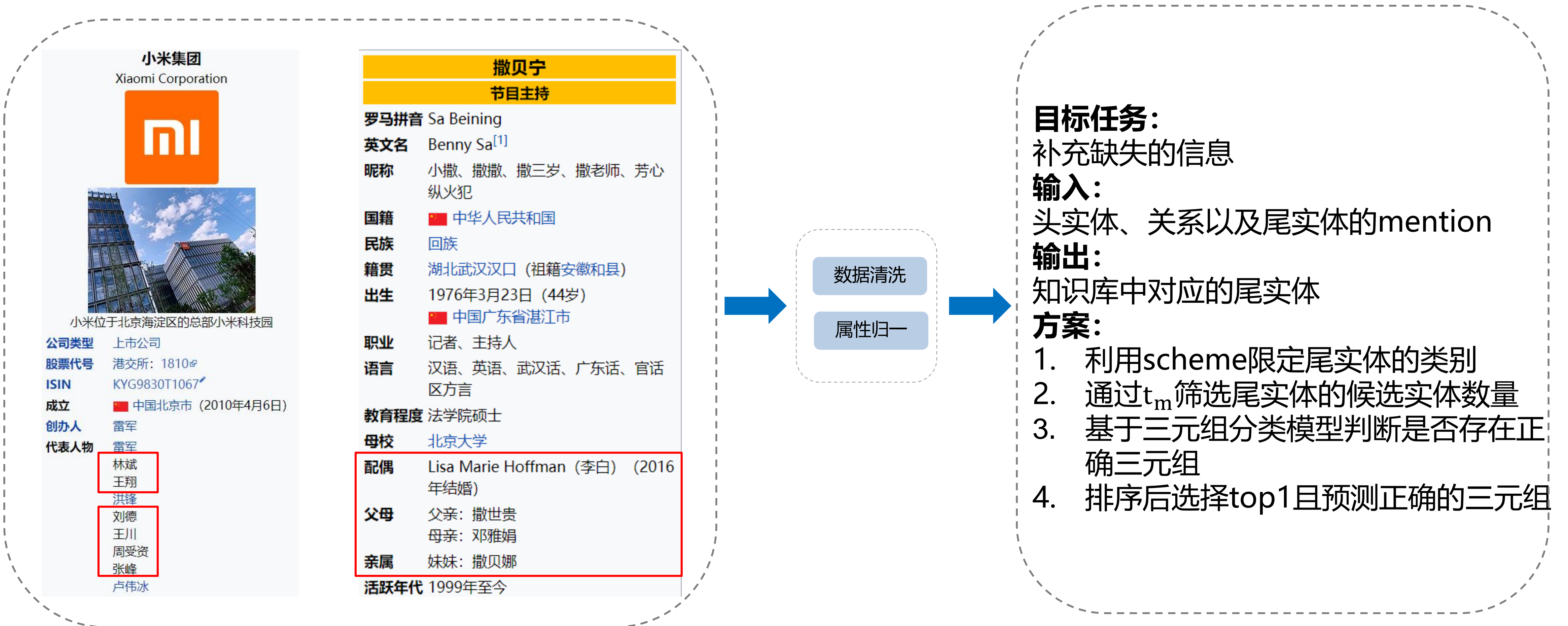
匹配模型



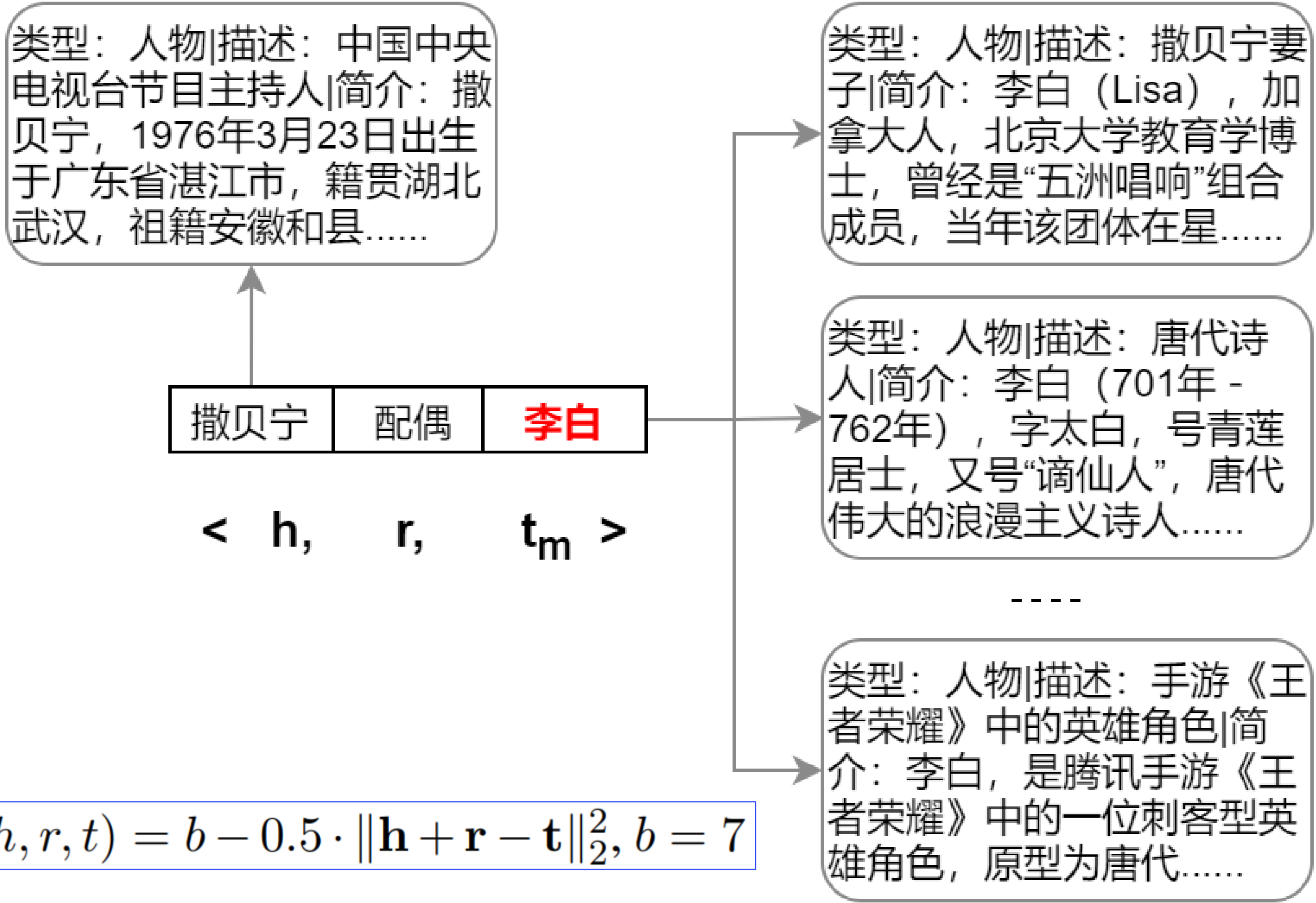
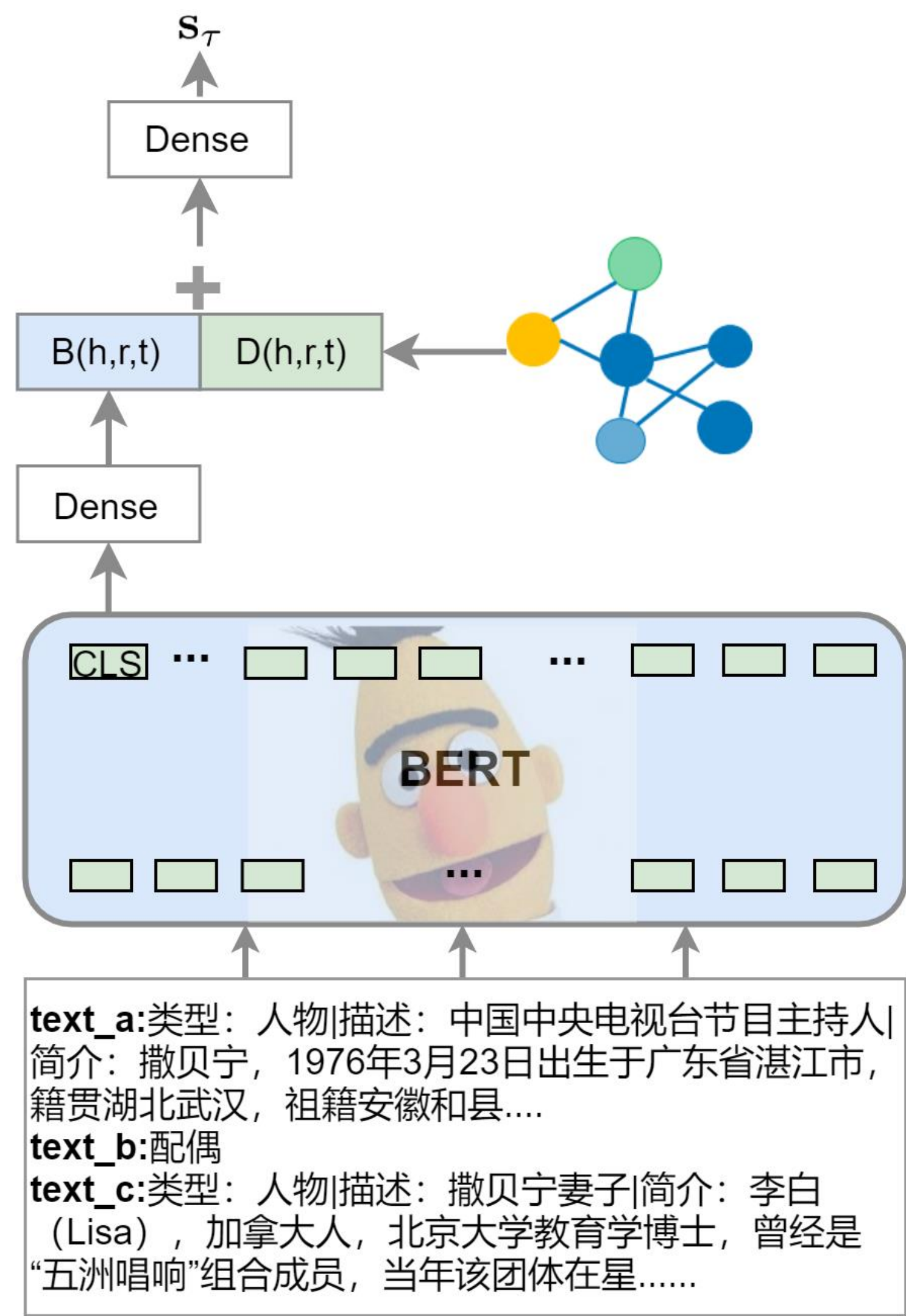
知识补全中的应用



知识图谱补全 (Knowledge Graph Completion, KGC) 目前主要被抽象成一个预测问题, 即预测出三元组中缺失的部分。



知识补全中的应用



$$z(h, r, t) = b - 0.5 \cdot \| \mathbf{h} + \mathbf{r} - \mathbf{t} \|_2^2, b = 7$$

$$D(h, r, t) = \frac{\exp\{z(h, r, t)\}}{\sum_{\bar{t} \in E_m} \exp\{z(h, r, \bar{t})\}}$$

$$\mathcal{L} = - \sum_{\tau \in \mathbb{D}^+ \cup \mathbb{D}^-} (y_\tau \log(s_{\tau 0}) + (1 - y_\tau) \log(s_{\tau 1}))$$

1	0.9	✓
0	0.2	
0	-	
0	0.1	

04

总结

总结与展望

1. 简单介绍了知识表示学习在实体链接、实体推荐和知识补全中的应用
2. 和word2vec一样，通过知识表示学习得到的实体向量可以应用到很多场景中
3. 工业界实用最重要，很多场景下，对模型的复杂性很敏感，选择模型需要综合考量
4. 知识表示学习的探索之路还有很长，同志们加油！！！！

1. Wang Z, Zhang J, Feng J, et al. Knowledge graph and text jointly embedding[C] //Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP). 2014: 1591-1601.
2. Zhong H, Zhang J, Wang Z, et al. Aligning knowledge and text embeddings by entity descriptions[C]//Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing. 2015: 267-272.
3. Xie R, Liu Z, Jia J, et al. Representation learning of knowledge graphs with entity descriptions[C]//Proceedings of the AAAI Conference on Artificial Intelligence. 2016, 30(1).
4. Xiao H, Huang M, Meng L, et al. SSP: semantic space projection for knowledge graph embedding with text descriptions[C]//Proceedings of the AAAI Conference on Artificial Intelligence. 2017, 31(1).
5. Reimers N, Gurevych I. Sentence-bert: Sentence embeddings using siamese bert-networks[J]. arXiv preprint arXiv:1908.10084, 2019.
6. Yao L, Mao C, Luo Y. KG-BERT: BERT for knowledge graph completion[J]. arXiv preprint arXiv:1909.03193, 2019.
7. 刘知远, 孙茂松, 林衍凯, 等. 知识表示学习研究进展[J]. 计算机研究与发展, 2016, 53(2): 247.



THANKS!

今天的分享就到这里...

Ending

