

Multimodal Knowledge Graphs: Construction, Inference, and Challenges

Meng Wang

Guilin Qi, Qiushuo Zheng, Chaoyu Bai

Southeast University

- **Multimodality**
- Multimodal KG Construction
- Inference
- Challenges

What is Multimodal Knowledge?

Multimodality: is the application of multiple literacies within one medium¹.



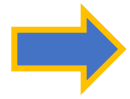
WIKIPEDIA
The Free Encyclopedia

Knowledge: **Facts** acquired through experience or education;
the theoretical or practical **understanding** of a **subject**

-----Oxford dictionary (English) , 2016

[1] <https://en.wikipedia.org/wiki/Multimodality>

**Multimodal knowledge:
is an awareness or understanding of someone
or something in different multimodalities.**



We can **extract different multimodal knowledge on a same fact (or a **same conventional knowledge**)**

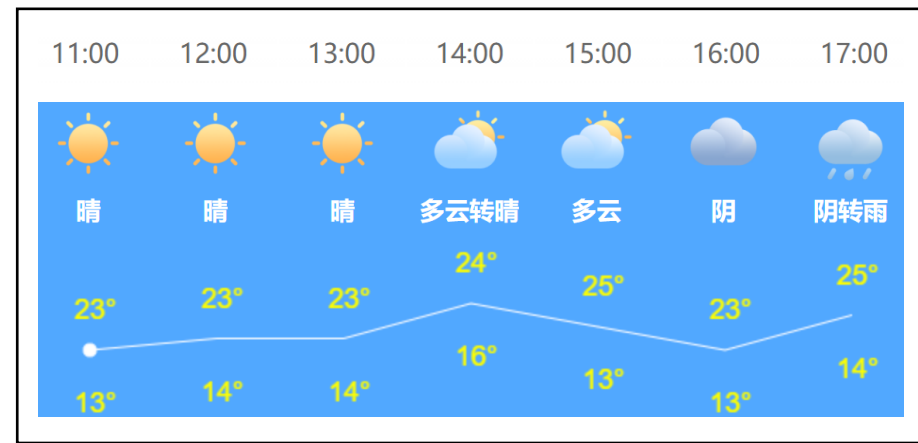
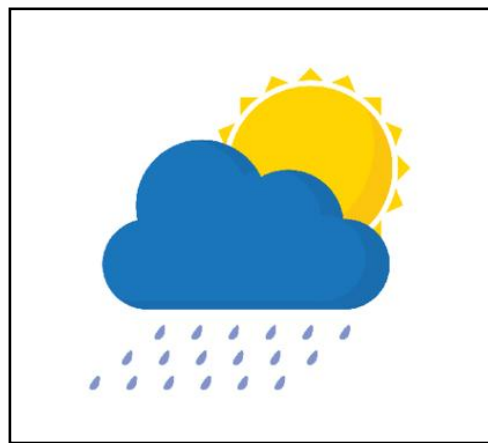
We can use different multimodal knowledge to understand on a same thing (a conventional entity)

Case 1:

南昌天气晴，局部地区有短时降雨



Spoken and Visual Knowledge

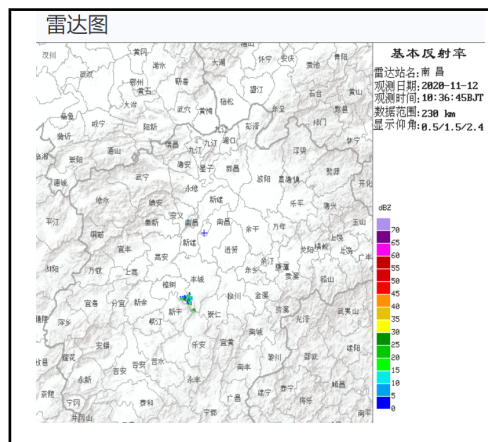


Symbol and Temperature Knowledge



Gestural and Linguistic Knowledge

欢迎大家参加
CCKS, 江西省
最近到秋天了,
所以南昌这几天
都是好天!



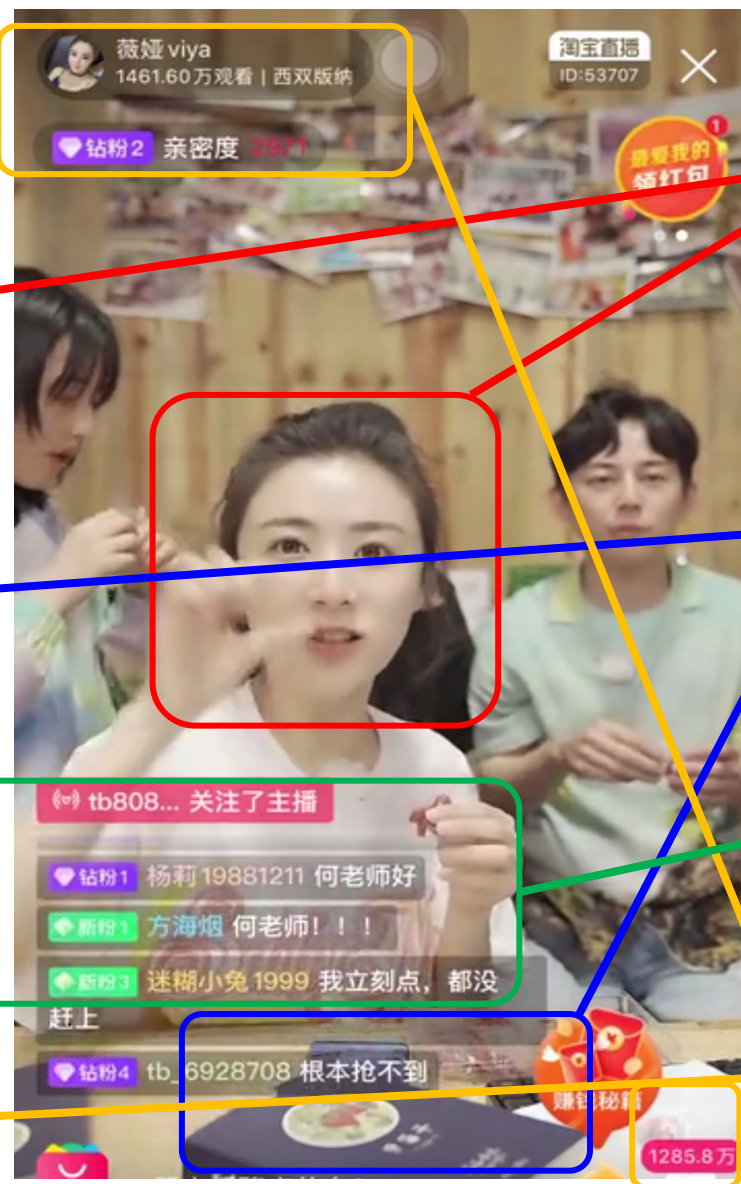
Geography Knowledge



Weather Specific Knowledge

(图片来源于网络, 仅供示意)

Case 2: 这款商品真的很不错！



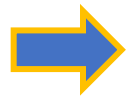
Gestural and Facial Expression Knowledge

Visual Knowledge

Text, Dialogue, and Linguistic Knowledge

Domain Specific Knowledge

We can extract different multimodal knowledge on a same fact (or a same conventional knowledge)



We can **use** different multimodal knowledge to understand on a same thing (**a conventional entity**)

Case 3:

London

Capital city



Clockwise from top: City of London in the foreground with Canary Wharf in the far background, Trafalgar Square, London Eye, Tower Bridge and a London Underground roundel in front of Elizabeth Tower



Location within Europe

- Show map of the United Kingdom
- Show map of Europe
- Show map of Earth

Coordinates: Show all  51°30'26"N 0°7'39"W

Sovereign State	United Kingdom
Constituent Country	England
Region	London (coterminous)
Counties	Greater London City of London
Settled by Romans	AD 47 ^[1] as <i>Londinium</i>
Districts	City of London & 32 boroughs
Government	
• Type	Executive mayoralty and deliberative assembly within unitary constitutional monarchy

Image

Text

KG

Case 4:

科比携手邓肯KG进入名人堂候选名单 三大传奇2020年将圆梦

2019
12/20
06:18

猫大熊
企鹅号

分享



评论

785



北京时间12月20日，奈史密斯篮球名人堂公布了2020年篮球名人堂候选人名单，其中包括科比、邓肯、加内特这三大NBA超级巨星，可谓是史上最星光熠熠的一届。



赛事精选 12/20 12/21 12/22

76人	未开始	独行侠	0	09:00	0
骑士	未开始	灰熊	0	08:00	0
步行者	未开始	国王	0	08:00	0
凯尔特人	未开始	活塞	0	08:30	0
猛龙	未开始	奇才	0	08:30	0

推荐视频

- 【球星】哈登集锦 53分16板17助历史第一人
- 【扣篮】位置感极佳 特纳底线埋伏好冲起补扣
- 【球星】庄神集锦 砍下25分18个篮板频频双手爆筐
- 【集锦】尼克斯 122-129火箭 哈登生涯新高 53+16+17破三分纪录
- 1日NBA头条 哈登 53+16+17历史第一 韦少半场三双

圣诞节大型立体粉色豪华圣诞树家用小型套装饰品摆件仿真1.5米

双旦礼遇季 此商品12.19开卖，请立即购买

价格 ¥428.00
礼遇价 ¥55.00 双旦礼遇季
淘金币可抵扣商品价格2%

运费 浙江金华 至 珠海 快递: 0.00

月销量 4.0万+ | 累计评价 54549 | 送天猫积分 27

颜色分类: [Image thumbnails]

数量 1 件 库存1147件

立即购买 加入购物车

服务承诺 正品保证 极速退款 赠运费险 支付方式

商品详情 累计评价 54549

品牌名称: 华驰

产品参数:

品牌: 华驰	尺寸: 141cm(含)-170cm(含)	颜色分类: 1米含52个配饰 1....
毛重: 2.2KG	光源类型: 发光	货号: SDS-001
包装体积: 0.04		

商品详情 累计评价 54549

与描述相符 太满意 便宜(447) 质量好(393) 快递不错(157) 态度不错(103) 做工不错(40) 颜色漂亮(32) 质量一般(34)

4.9

全部 追评 (934) 图片 (9343)

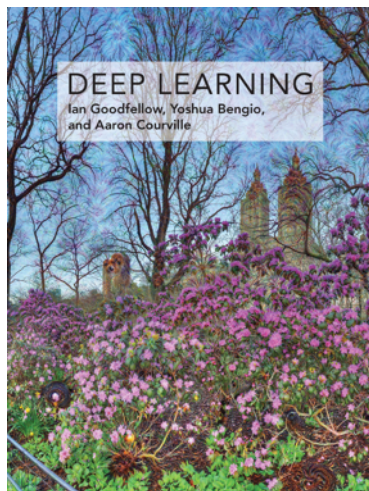
女儿想要节日的气氛就给买了个，虽然安装的时候树枝有点扎手，但装好后确实漂亮，尤其晚上把彩灯打开后超级美，女女非常喜欢。

颜色分类: 1.5米 h***1 (图)

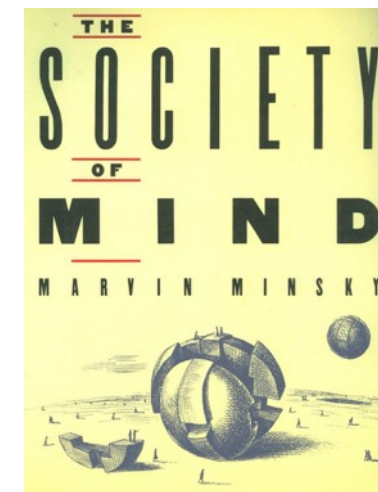
解释: 感谢您对我们的支持与信赖，您简短的几句话胜过无数的赞美，您的满意就是对我们最大的支持！希望您带去的不仅仅是一份舒适分感觉，更是一份美丽的心情！

Why we need multimodality?

Cognitive and Knowledge Graph View



Yoshua Bengio
NeurIPS Keynote, 2019



Marvin Minsky
The Society of Mind, 1986

SYSTEM 1 VS. SYSTEM 2 COGNITION

2 systems (and categories of cognitive tasks):

System 1

- Intuitive, fast, **UNCONSCIOUS**, non-linguistic, habitual
- Current DL

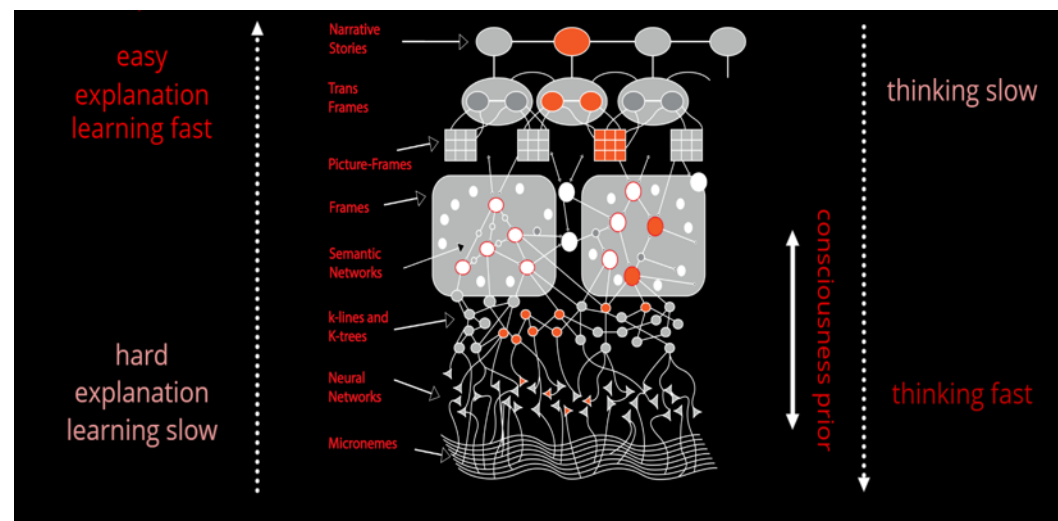
System 2

- Slow, logical, sequential, **CONSCIOUS**, linguistic, algorithmic, planning, reasoning
- Future DL

Manipulates high-level / semantic concepts, which can be recombined combinatorially

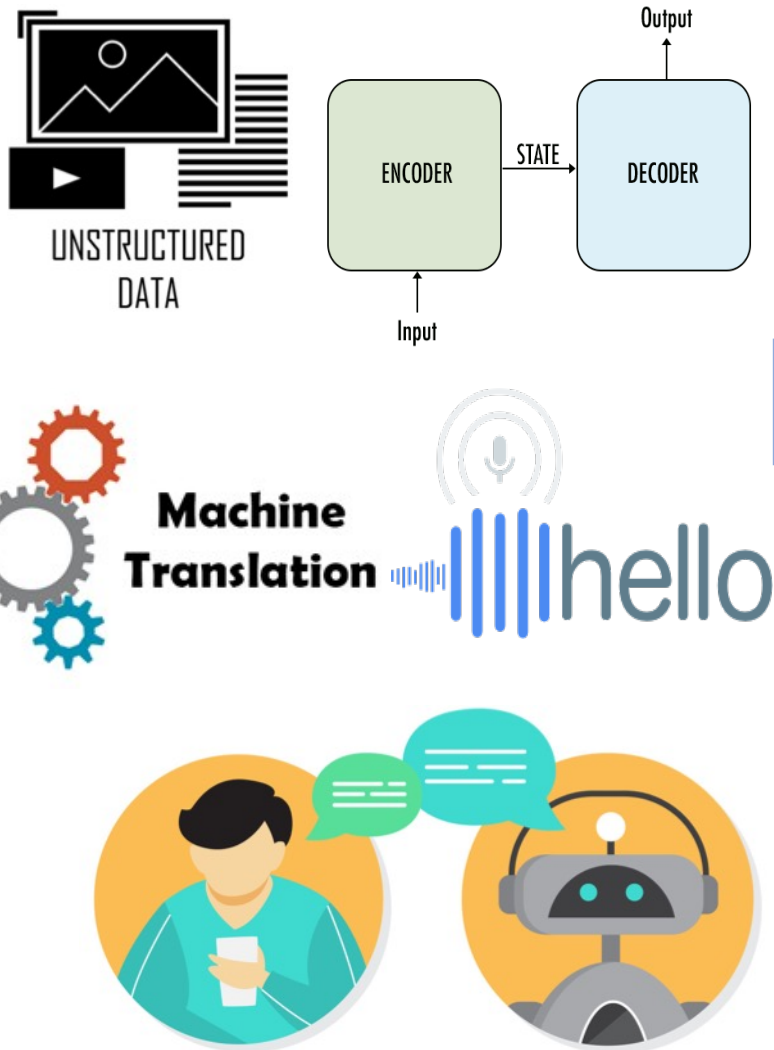
Mila

From system 1 DL to system 2 DL



Framework for representing knowledge

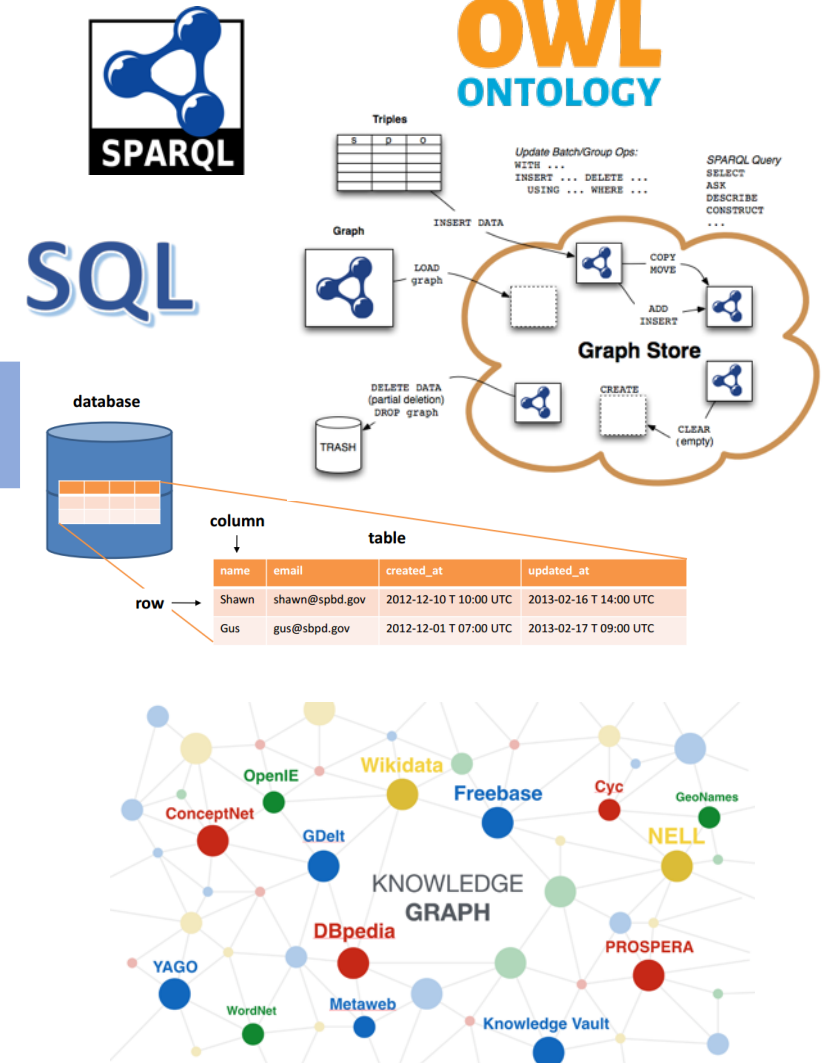
Neural (System1)

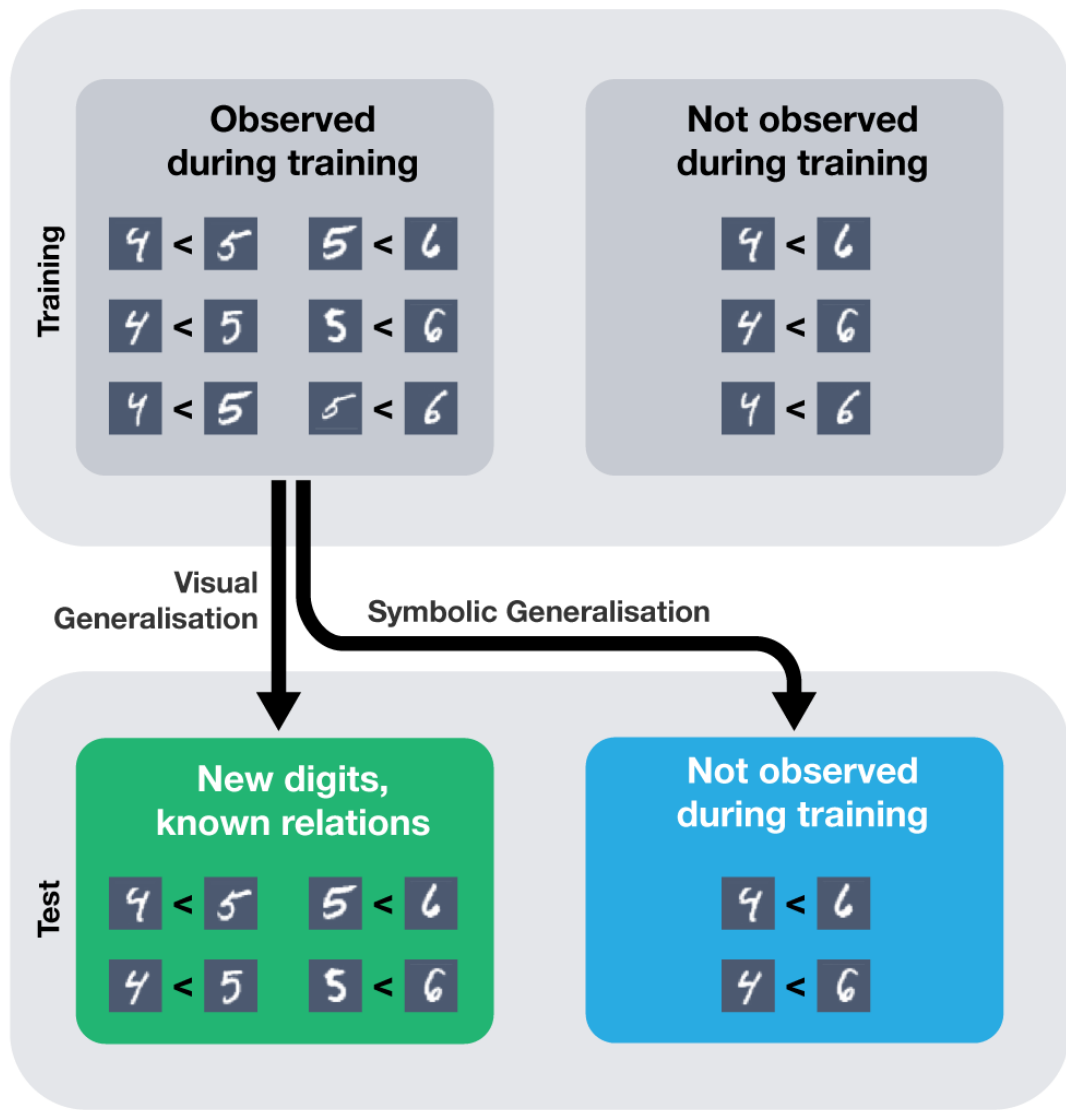


Input
(Question)

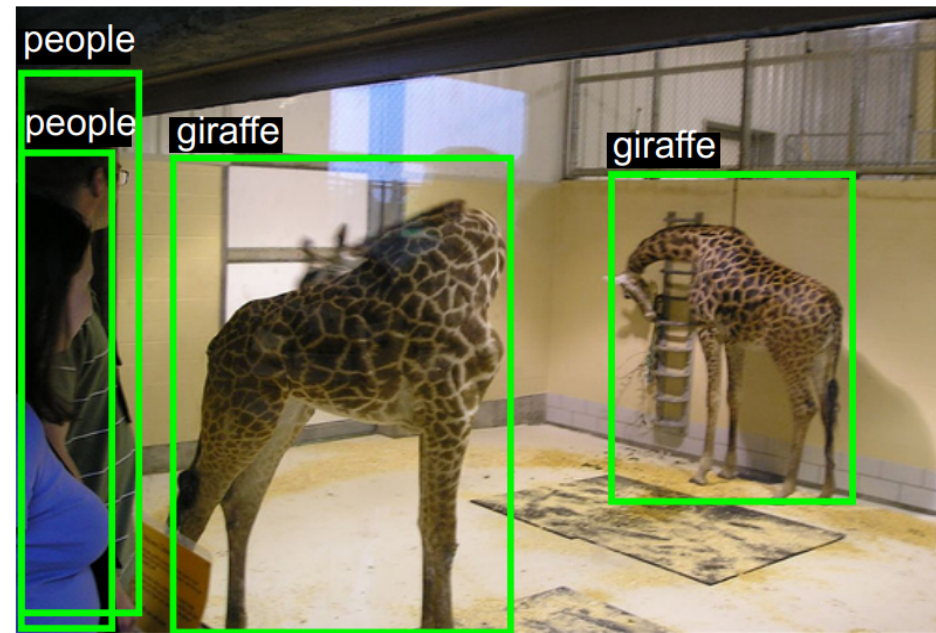
Goal
(Answer)

Symbolic (System2)





Visual generalisation vs. Symbolic generalisation



Attributes:

- glass
- house
- room
- standing
- walking
- wall
- zoo

Scenes:

- museum
- indoor

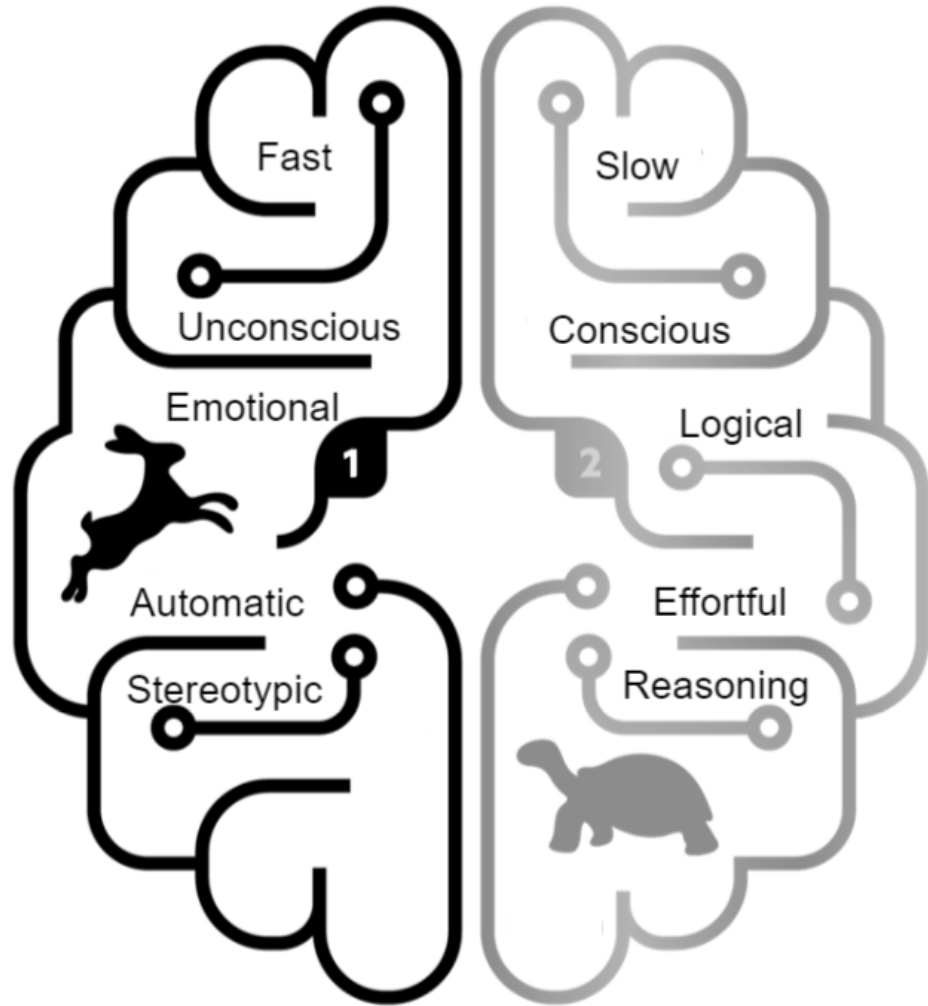
Visual Question: How many giraffes are there in the image?
Answer: Two.

Common-Sense Question: Is this image related to zoology?
Answer: Yes. **Reason:** Object/Giraffe --> Herbivorous animals --> Animal --> Zoology; Attribute/Zoo --> Zoology.

KB-Knowledge Question: What are the common properties between the animal in this image and zebra?
Answer: Herbivorous animals; Animals; Megafauna of Africa.

VQA, Commonsense QA, KBQA, and Machine Reading Comprehension

Cognitive Theory



Knowledge Graph Perspective

Neural (system1) are

- powerful for some problems
- robust to data noise
- hard to understand or explain
- poor at symbol manipulation
- unclear how to effectively use background knowledge

Symbolic (system2) are

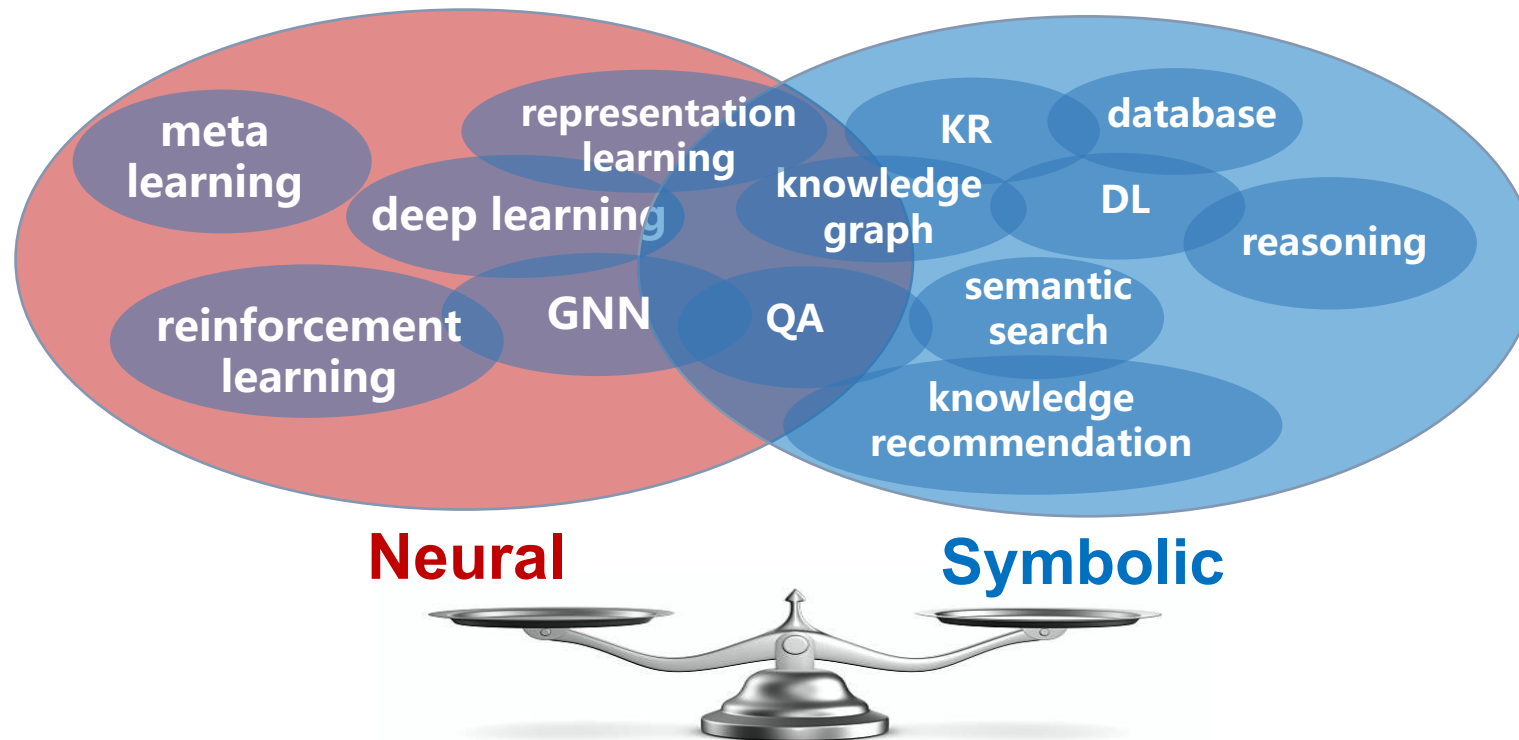
- Usually poor regarding machine learning problems
- Intolerant to data noise
- Easy to understand and assess by a human
- Good at symbol manipulation
- Designed to work with background knowledge

Neural+Symbolic:

- powerful machine learning paradigm
- robust to data noise
- easy to understand and assess by humans
- good at symbol manipulation
- work seamlessly with background knowledge

HOW TO

Multimodal Knowledge Graph ?



Application View

**More cross-modal relations, more details
and more answers**

London

Capital city



122 Leadenhall Street

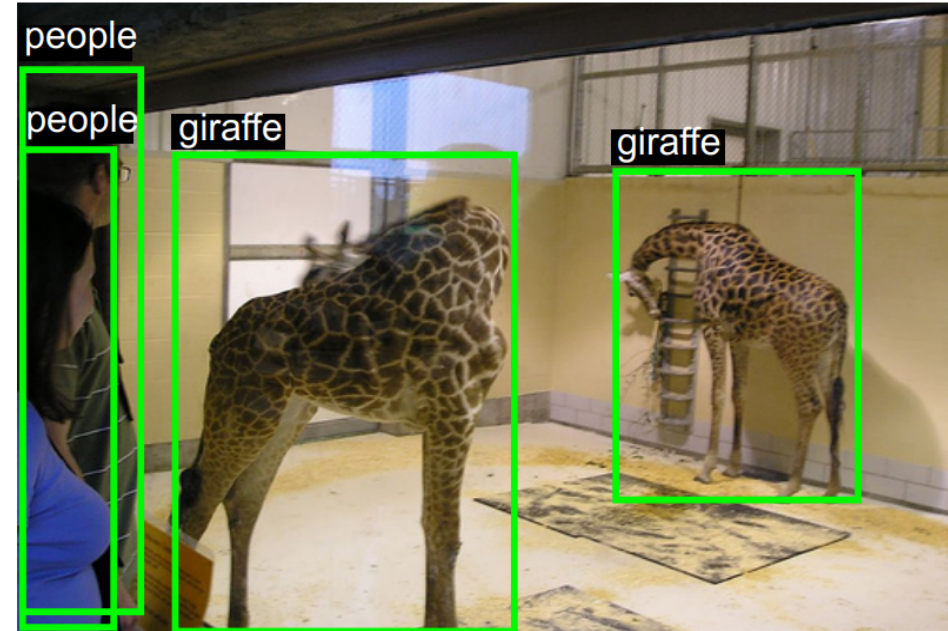
From Wikipedia, the free encyclopedia
(Redirected from [Leadenhall Building](#))

122 Leadenhall Street, also known as the [Leadenhall Building](#), is a [skyscraper](#) i opened in July 2014 and was designed by [Rogers Stirk Harbour + Partners](#); it is because of its distinctive wedge shape similar to that of [the kitchen utensil](#) with tall buildings recently completed or under construction in the [City of London](#) fin [Street](#), [The Pinnacle](#), and [The Scalpel](#).

122 Leadenhall Street



Cross-modal entity grounding



Attributes:

glass
house
room
standing
walking
wall
zoo

Scenes:

museum
indoor

Visual Question: How many giraffes are there in the image?

Answer: Two.

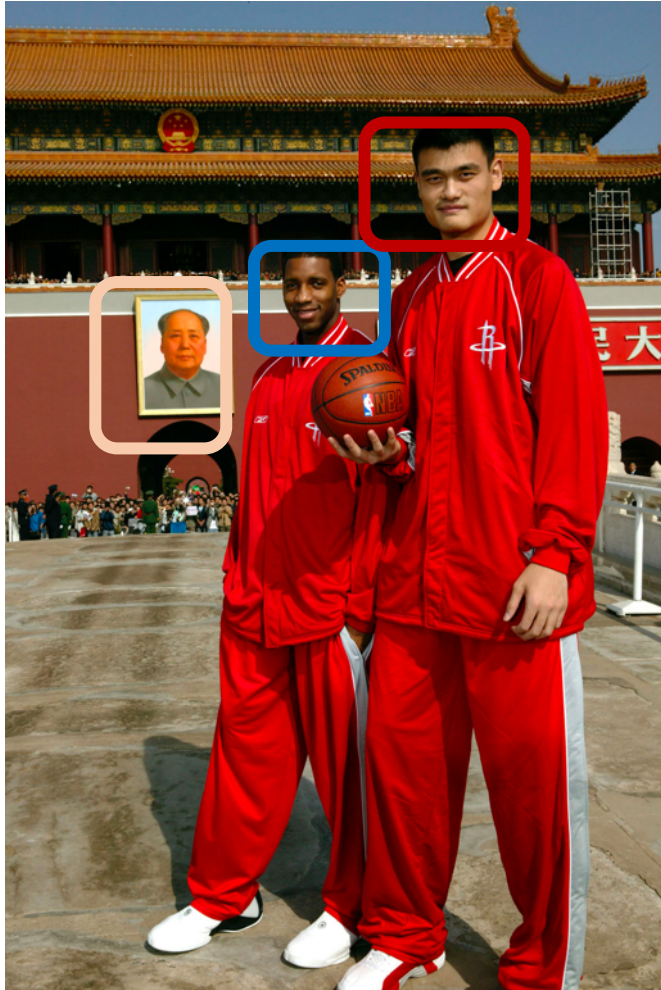
Common-Sense Question: Is this image related to zoology?

Answer: Yes. **Reason:** Object/Giraffe --> Herbivorous animals --> Animal --> Zoology; Attribute/Zoo --> Zoology.

KB-Knowledge Question: What are the common properties between the animal in this image and zebra?

Answer: Herbivorous animals; Animals; Megafauna of Africa.

VQA, Complicated scene understanding



Tracy McGrady
Basketball Player

Tracy Lamar McGrady Jr. is an American former professional basketball player. He is best known for his career in the National Basketball Association, where he played as both a shooting guard and small forward. McGrady is a seven-time NBA All-Star, seven-time All-NBA selection, two-time NBA scoring champion, and on...



Yao Ming

Yao Ming is a Chinese basketball executive and retired professional basketball player who played for the Shanghai Sharks of the Chinese Basketball Association and the Houston Rockets of the National Basketball Association. He was selected to start for the Western Conference in the NBA All-Star Game eight times, a...

Wikipedia Twitter Facebook



\$645.00
Stella McCartney Loop L...
Neiman Marcus



\$22.33
Sport-tek Men's Elastic D...



Yao Ming and Tracy McGrady of the Houston Rockets visit Beijing in 2004

All

Looks like

Similar images

Pages with this

Related searches

Looks like



Yao Ming



Tracy McGrady
Basketball Player



Mao Zedong
Former President of the People's Republic of C...

See more

Similar images



Remove Crop



Remove Crop

All

Shop for similar

Similar images

Related searches



\$645.00
Stella McCartney Loop L...
Neiman Marcus



\$442.00
Jil Sander Leather Platfo...
Lyst



\$60.00
Infant Vans Classic Slip-...
Walmart



\$81.52
Lacoste Men's Marice Bl ...
2daydeliver.com



\$39.99
Crocs Pfd White Speciali...
crocs.com



\$45.00
Women's Keds Champio...
shoes.com



**Cross-modal Disambiguation:
Heterogeneous in modal, but correlated in semantic**



API 中心

搜索相关内容

简介

API 概览

调用方式

图像处理相关接口

图像审核相关接口

图像理解相关接口

· 公众人物识别

· 图像标签

文档中心 > API 中心 > 图像分析 > 图像理解相关接口 > 公众人物识别

公众人物识别

最近更新时间: 2019-08-22 19:41:50

1. 接口描述

接口请求域名: tiia.tencentcloudapi.com。

传入一张图片, 可以识别图片中包含的人物是否为公众人物, 如果是, 输出人物的姓名, 支持识别一张图片中存在的多个人脸, 针对每个人脸, 会给出与之最相似的公众人物。

默认接口请求频率限制: 20次/秒。

Result :

姚明 (100)

威尔史密斯 (39)

梅兰芳 (36)

威尔·史密斯

演员



Visual Entity Disambiguation



刘欢在美国超市被偶遇，买8美元面包，爽快接过纸笔给网友签名

Textual Entity Disambiguation

刘欢 +

这是一个多义词，请在下列义项上选择浏览（共14个义项）

- 刘欢：中国内地流行音乐家
- 刘欢：广东省广州市中级人民法院助理审判员
- 刘欢：中国足球运动员
- 刘欢：长虹街道办事处副主任
- 刘欢：湖南发展研究中心研究员联络处副主任
- 刘欢：清华大学环境学院副教授
- 刘欢：清华大学教师
- 刘欢：象棋棋手
- 刘欢：矿大（北京）管院第十二届研究生会副主席
- 刘欢：苏州东吴队球员
- 刘欢：中国大陆男演员
- 刘欢：扣篮王刘欢
- 刘欢：全国技术能手



基本信息

中文名	刘欢	毕业院校	国际关系学院法国文学专业
外文名	Liu Huan	经纪公司	百娱传媒股份有限公司
别名	欢哥	代表作品	少年壮志不言愁、弯弯的月亮、心中的太阳、千万次的问、这一拜、好汉歌、从头再来、凤凰于飞
国籍	中国	主要成就	CCTV MTV音乐盛典最受欢迎男歌手
民族	汉族		《音乐风云榜》终身成就奖
星座	处女座		北艺协会电视剧优秀音乐创作奖
血型	O型		第十届华语歌曲“榜中榜”之“评委会特别奖”
身高	173cm		第四届中国金唱片“最佳流行专辑”
出生地	天津	生肖	兔
出生日期	1963年8月26日		
职业	歌唱家、音乐人、词曲创作人、大学音乐教授		

Huan Liu

Computer scientist

Huan Liu is a computer scientist at Arizona State University in Tempe, Arizona. He was named a Fellow of the Institute of Electrical and Electronics Engineers in 2012 for his contributions to feature selection in data mining and knowledge discovery. [Wikipedia](#)

What is Multimodal Knowledge Graph?

Richpedia

Still many unsolved problems

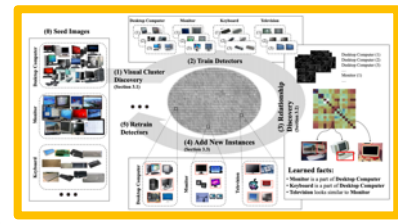
Multimodal Knowledge Graph 2019

IMGpedia 2015,2017

Semantic Web formats



NEIL: Image Knowledge Miner 2013



- Automatically extracting
- Semi-supervised learning
- Discovers common sense relationships

ImageNet, Visipedia 2009,2010

built upon the backbone of the WordNet



- Interlinking Multimedia
- Apply Linked Data Principles to Multimedia Fragments

iM 2008, 2009

ESP Game 2004

- Labeling Images with a Computer Game



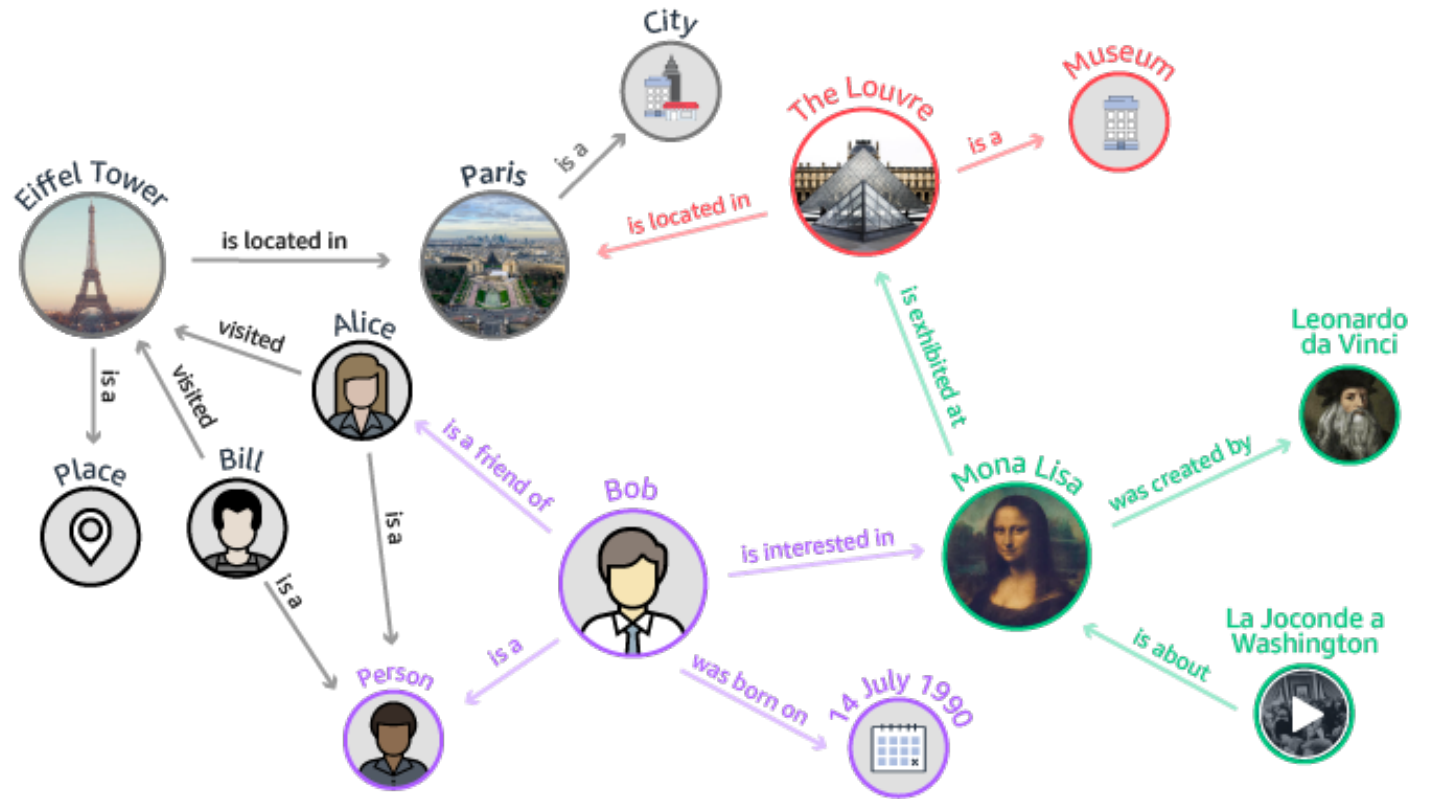
Emerging Multimodal KG Work

Node:

- Image entity
- Text entity
- Visual concept
- Textual concept

Relation:

- is-a
- has-visual-object
- meta-of
- has-tag
- co-locate-with



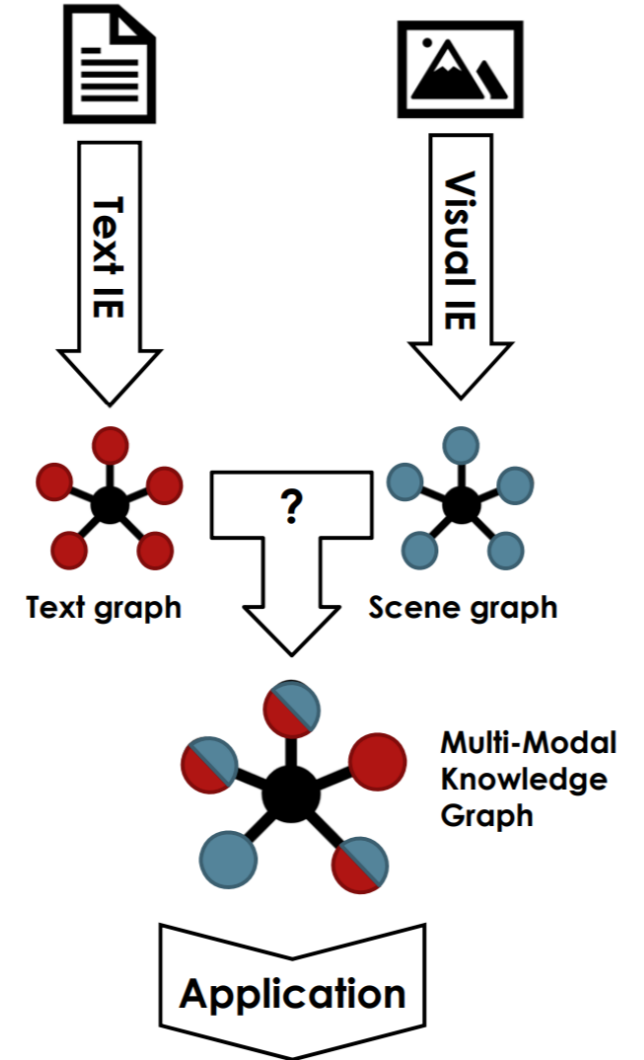
Emerging Multimodal KG Work

▶ Challenges:

- ▶ Parsing text to structured semantic graph
- ▶ Parsing images/videos to structures
- ▶ Grounding event/entities across modalities
- ▶ Multimodal argument role

▶ Applications

- ▶ Story Generation and Summarization
- ▶ Question Answering
- ▶ Commonsense Discovery



Emerging Multimodal KG Work

Home Back Query: Target: Януковича (Yanukovych) Event Search Number of Events: 2

Automated Summary: Source Document Translation from Ukrainian/Russian Show Visual Knowledge Elements Hide Visual Knowledge Elements

Event Summary: Столкновения 20 февраля стали одним из ключевых факторов, вынудивших Президента Украины Виктора Януковича пойти на подписание Соглашения об урегулировании политического кризиса на Украине, потере доверия к самому Януковичу и к реформатированию парламентского большинства. 20 февраля постановление о запрете применения силы властью (Translation: The clashes on February 20 became one of the key factors that forced President of Ukraine Viktor Yanukovich to sign the Agreement on the settlement of the political crisis in Ukraine, the loss of confidence in Yanukovich himself and the reformatting of the parliamentary majority, which issued a resolution on the evening of February 20 banning the use of force), что согласно всем имеющимся уликам те милиционеры и демонстранты, что стали жертвами снайперского огня, застрелены одними и теми же снайперами (Translation: that according to all available evidence, those policemen and demonstrators who became victims of sniper fire were shot by the same snipers)

Visual Entity Linking: Visual Entity Extraction: Event Arguments: Event Type

Source Doc & Text Extraction Result: Source Doc Translation: Recommended Events

Knowledge Elements based Ranking Incorporating User Feedback Similar Events Dissimilar Events

Date Location country
 Attackers his
 Target Unknown
 Instrument Unknown
 Type of Attack Conflict.Attack

HC000T6CP, 2011-01-19
 Event Time Person Organization GeopoliticalEntity Location
 Facility Vehicle Weapon Other

1 The case of Kiev snipers
 2 Red Cross Volunteers of Ukraine provide first aid to a wounded man on Institska the beginning of the eleventh hour on February 20, 2014
 3 Self-defenders carry out comrade on Institskaya Street to the rear at the end of eleventh hour on February 20
 4 Mark 13 on a pierced bullet on Institskaya Street, pasted by criminologists near the side opposite the Maidan
 5 The case of Kiev snipers question about the organizers and perpetrators of sniper: Euromaidan participants and at the same time law enforcement officers in Kiev on 20, 2014, which killed 53 people (49 protesters and 4 law enforcement officers)

Instrument Unknown
 Type of Attack Conflict.Attack

Date 201402
 Location Unknown
 Attackers group
 Target Unknown
 Instrument Unknown
 Type of Attack Conflict.Attack

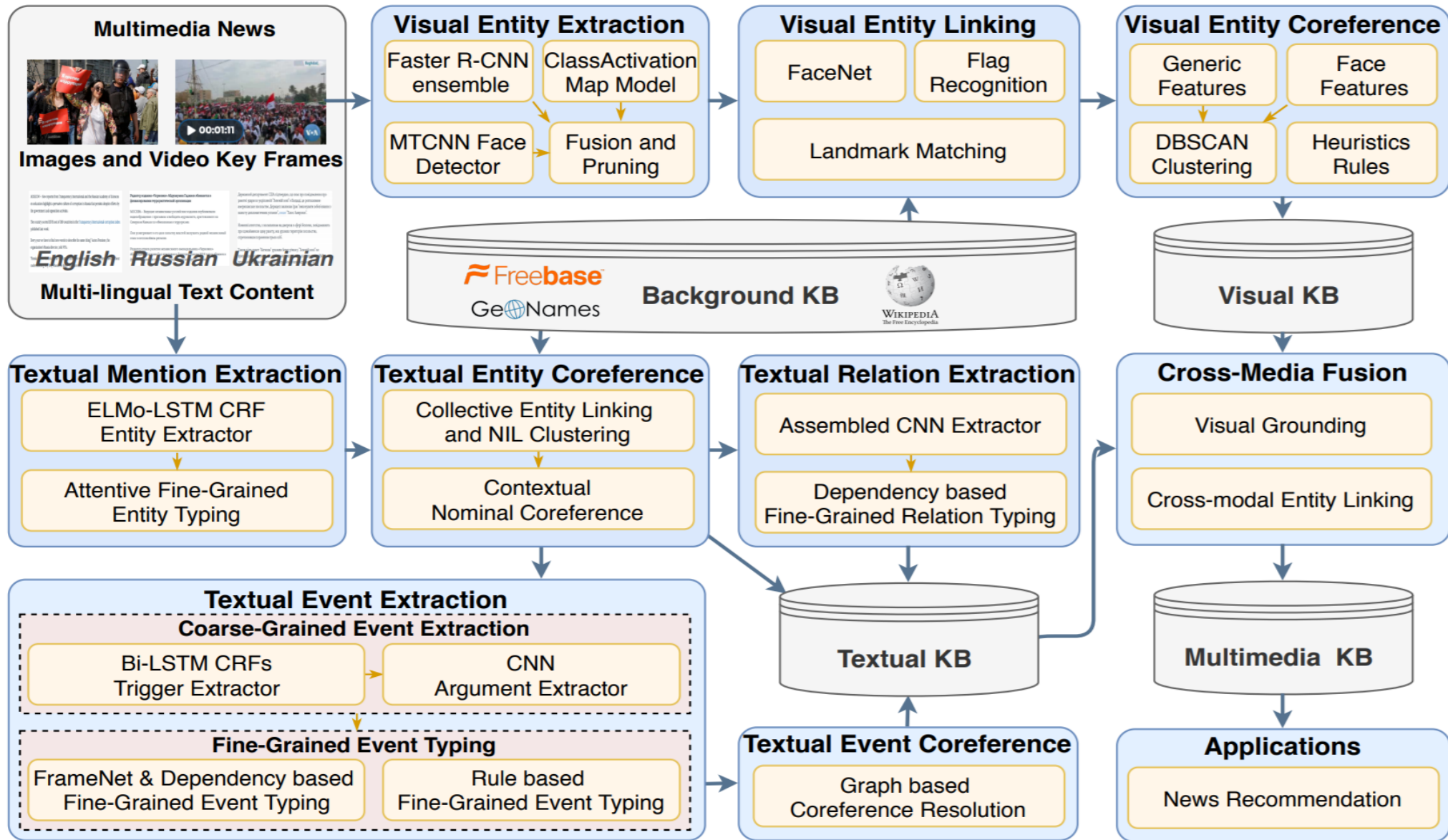
Date Location Attackers Target Instrument Type of Attack
 201402 Unknown Unknown Януковича (Yanukovych) Unknown Conflict.Attack

	Coarse-grained Types	Fine-grained Types
Entity	7	187
Relation	23	61
Event	47	144



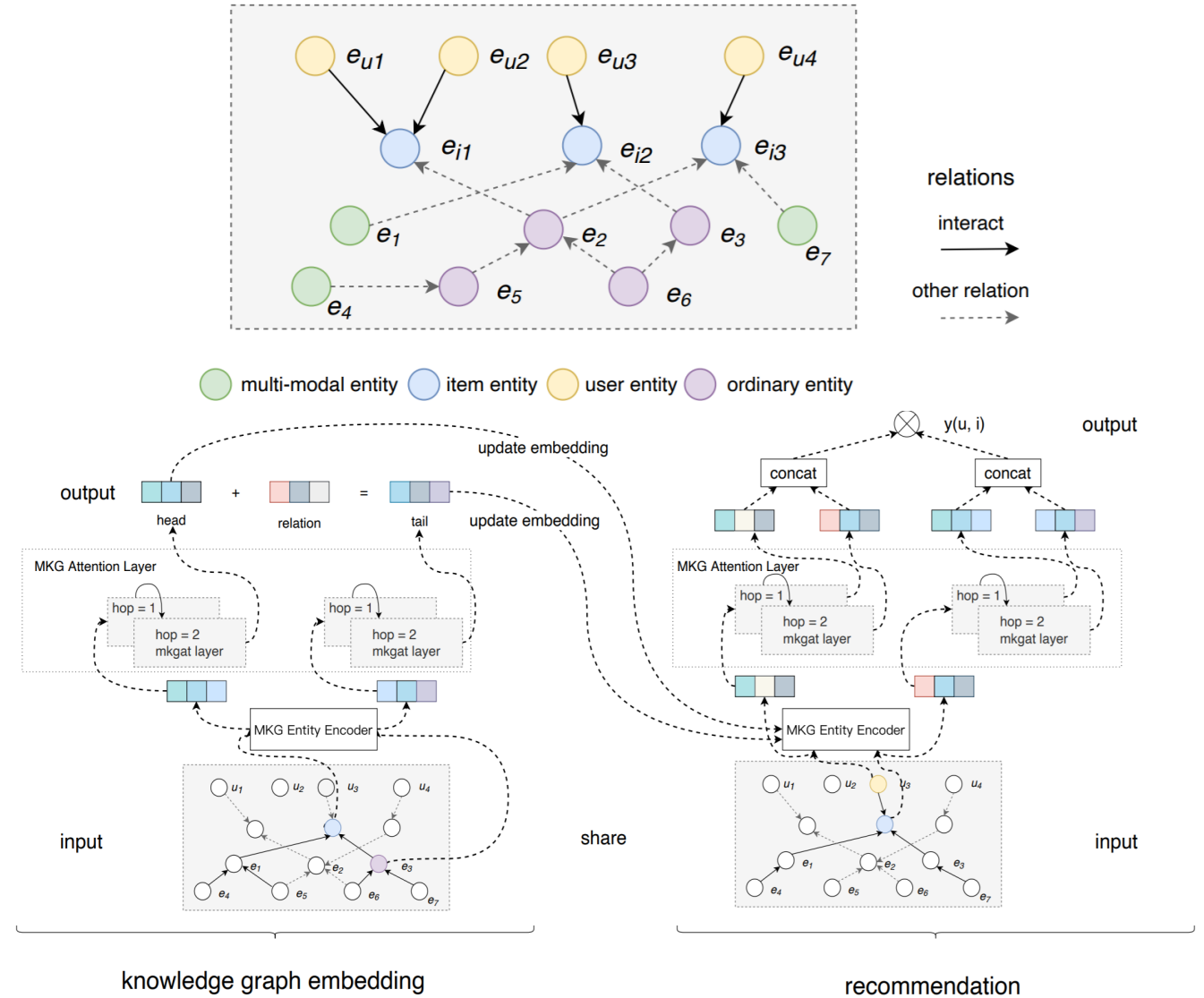
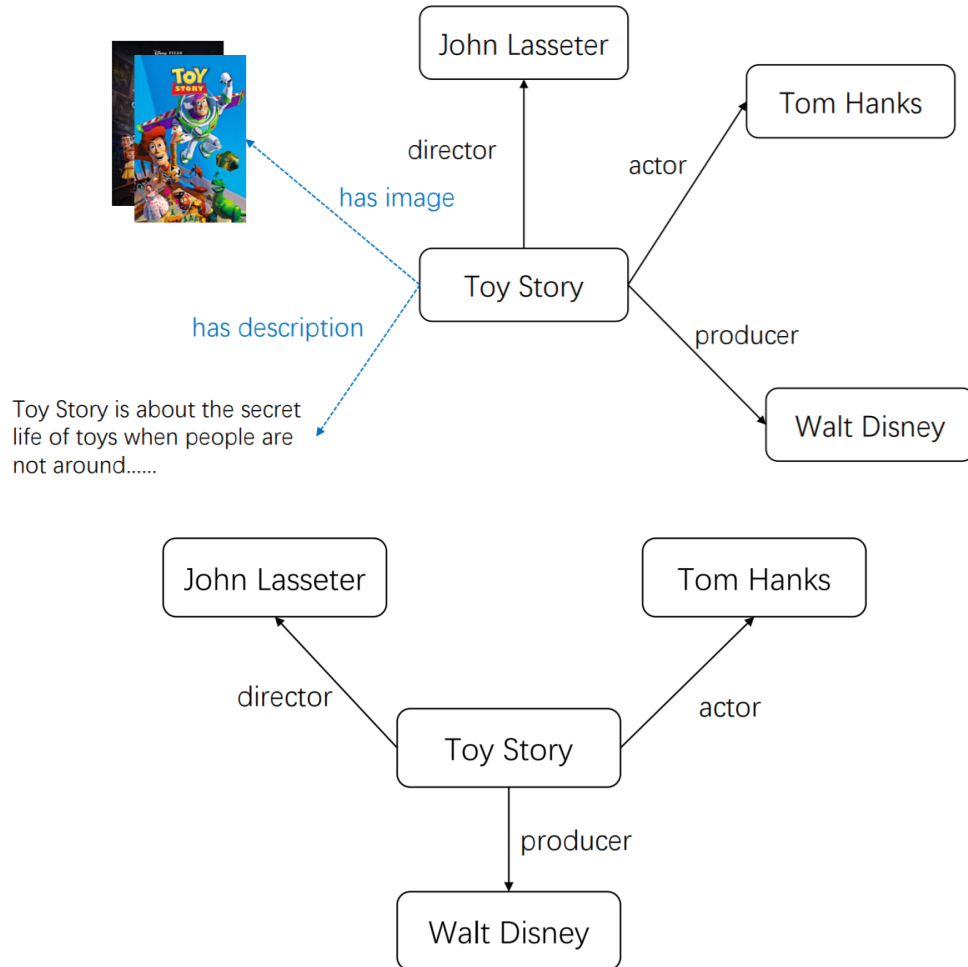
Li, Manling, et al. "Gaia: A fine-grained multimedia knowledge extraction system." *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*. (ACL 2020).

Emerging Multimodal KG Work



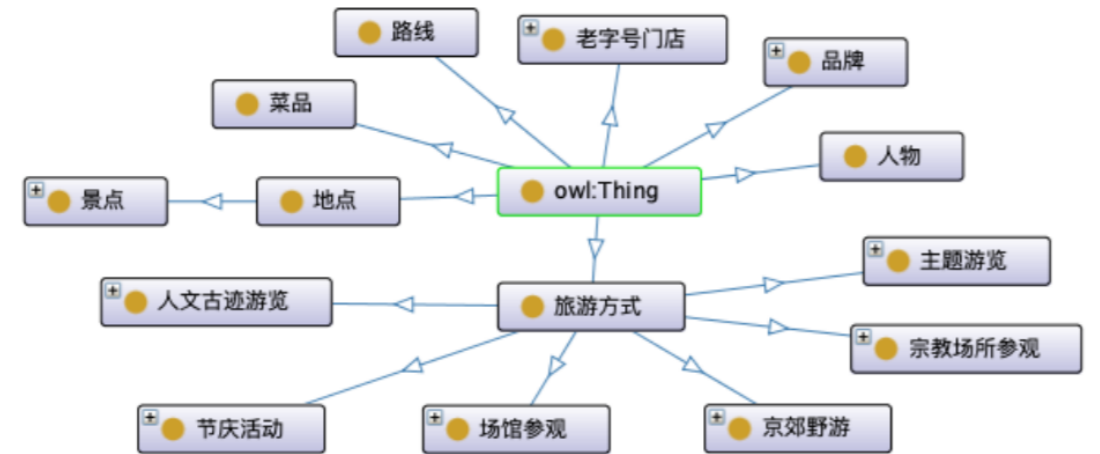
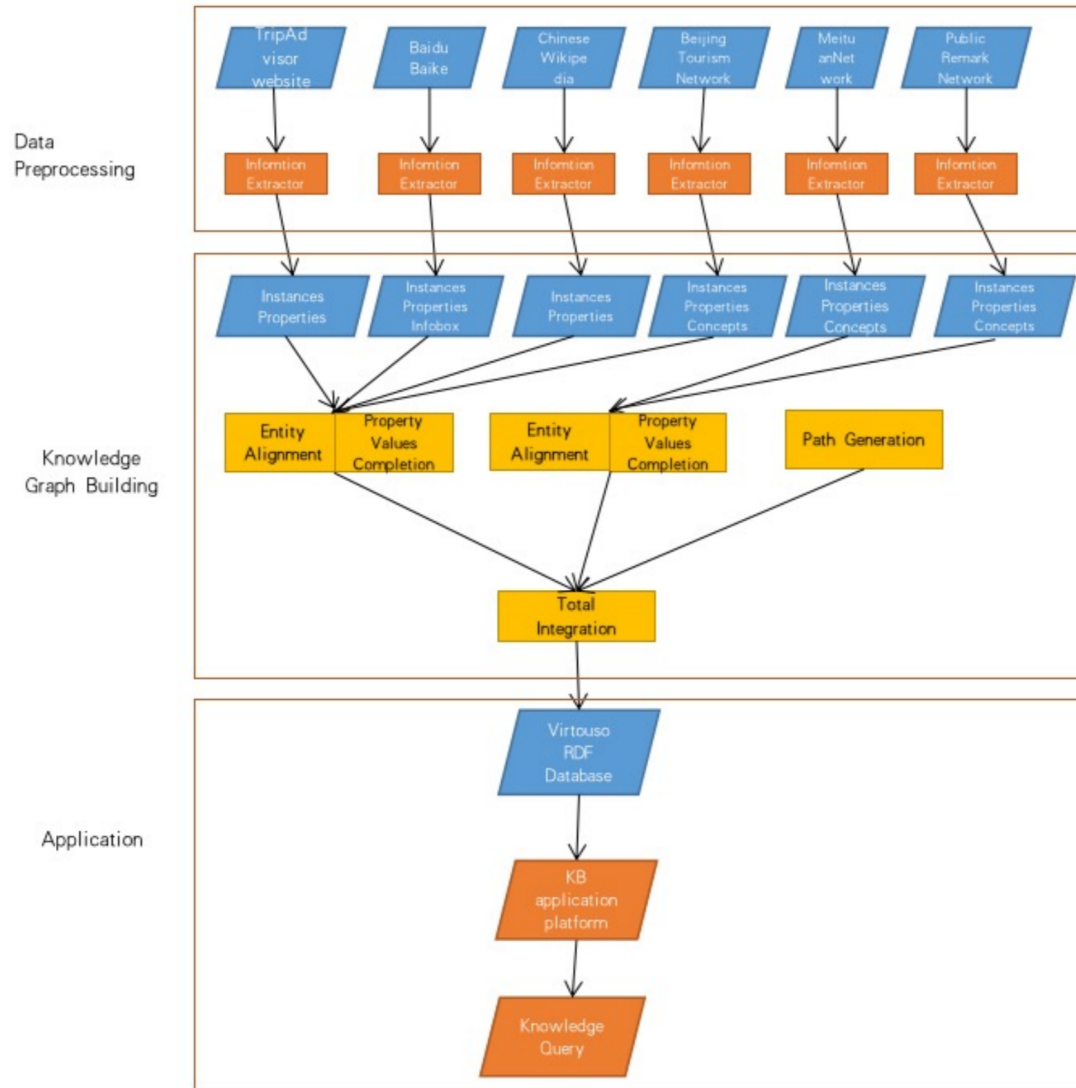
Li, Manling, et al. "Gaia: A fine-grained multimedia knowledge extraction system." *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*. (ACL 2020).

Emerging Multimodal KG Work



Sun, Rui, et al. "Multi-modal Knowledge Graphs for Recommender Systems." *Proceedings of the 29th ACM International Conference on Information & Knowledge Management (CIKM 2020)*.

Emerging Multimodal KG Work



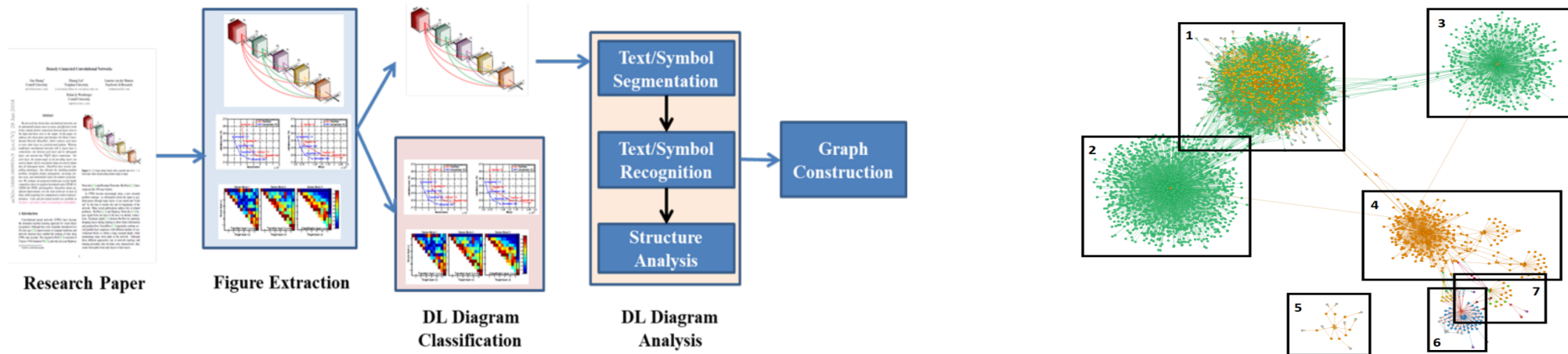
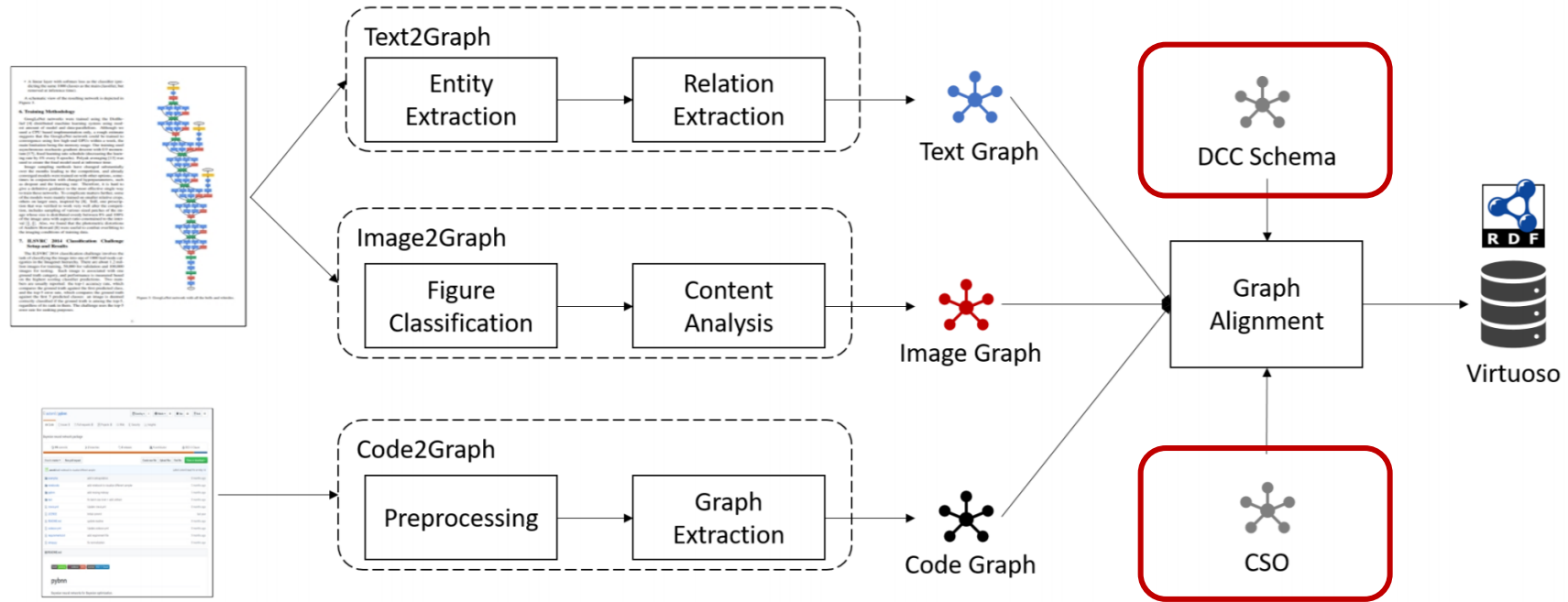
景点示例

<p>故宫博物院</p>  <p>概念类型: 博物馆世界 文化遗址 地理位置: 北京市东城区景山前街4号 联系电话: 010-65132255 景点级别: AAAAA级</p>	<p>圆明园公园</p>  <p>概念类型: 公园文化 节庆场所 皇家园林 地理位置: 清华西路28号 联系电话: 010-82670330 景点级别: AAAA级</p>	<p>玉渊潭公园</p>  <p>概念类型: 亲子游景区 文化节庆场所 野外景色 地理位置: 北京市海淀区西三环中路10号(中央电视台塔对面) 联系电话: 010-88653775, 010-88653806 景点级别: AAAA级</p>
--	--	--

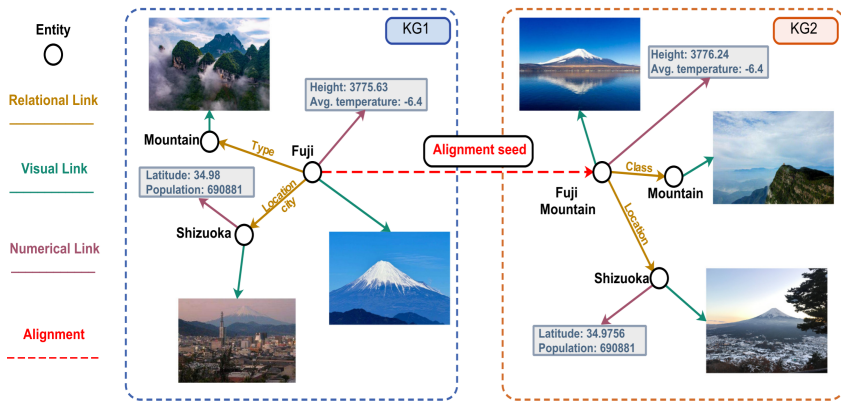
发展情况



Emerging Multimodal KG Work



Emerging Multimodal KG Work

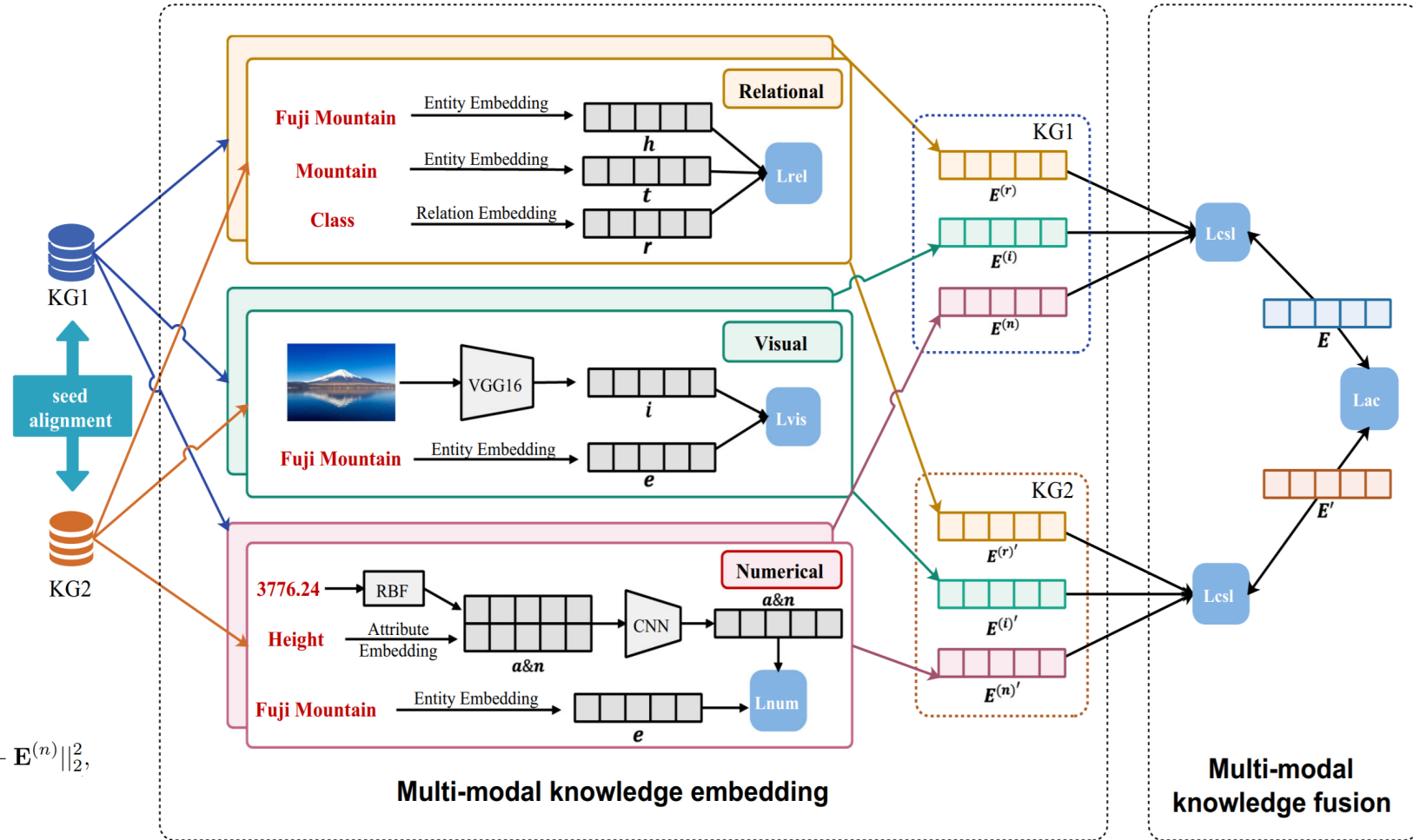


Datasets:
MMKG FB15K-DB15K and
FB15K-YAGO15K

$$\mathcal{L} = - \sum_{t \in \mathbf{T}} \log p(t \mid \theta_1, \dots, \theta_n).$$

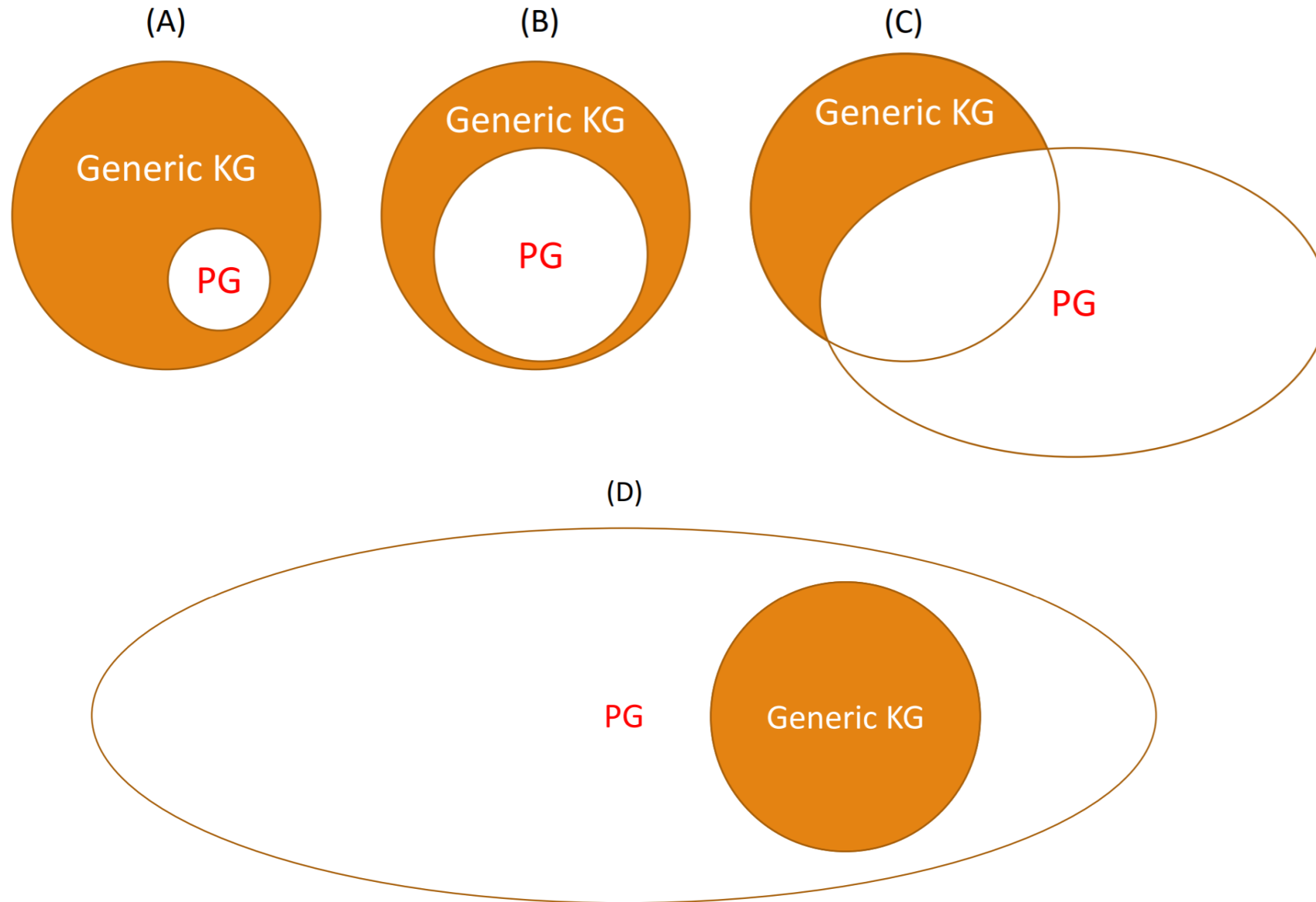


$$L_{csl}(\mathbf{E}, \mathbf{E}^{(r)}, \mathbf{E}^{(i)}, \mathbf{E}^{(n)}) = \alpha_1 \|\mathbf{E} - \mathbf{E}^{(r)}\|_2^2 + \alpha_2 \|\mathbf{E} - \mathbf{E}^{(i)}\|_2^2 + \alpha_3 \|\mathbf{E} - \mathbf{E}^{(n)}\|_2^2,$$



Liu, Ye, et al. “MMKG: multi-modal knowledge graphs.” *European Semantic Web Conference (ESWC 2019)*.
Chen, Liyi, et al. “MMEA: Entity Alignment for Multi-modal Knowledge Graph.” *International Conference on Knowledge Science, Engineering and Management (KSEM 2020)*. **(Best Paper)**

Emerging Multimodal KG Work



Emerging Multimodal KG Work

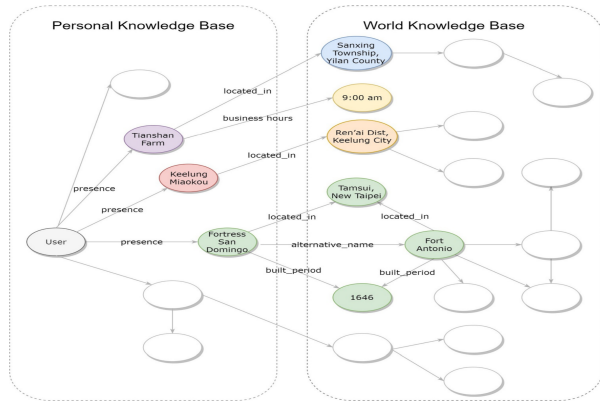
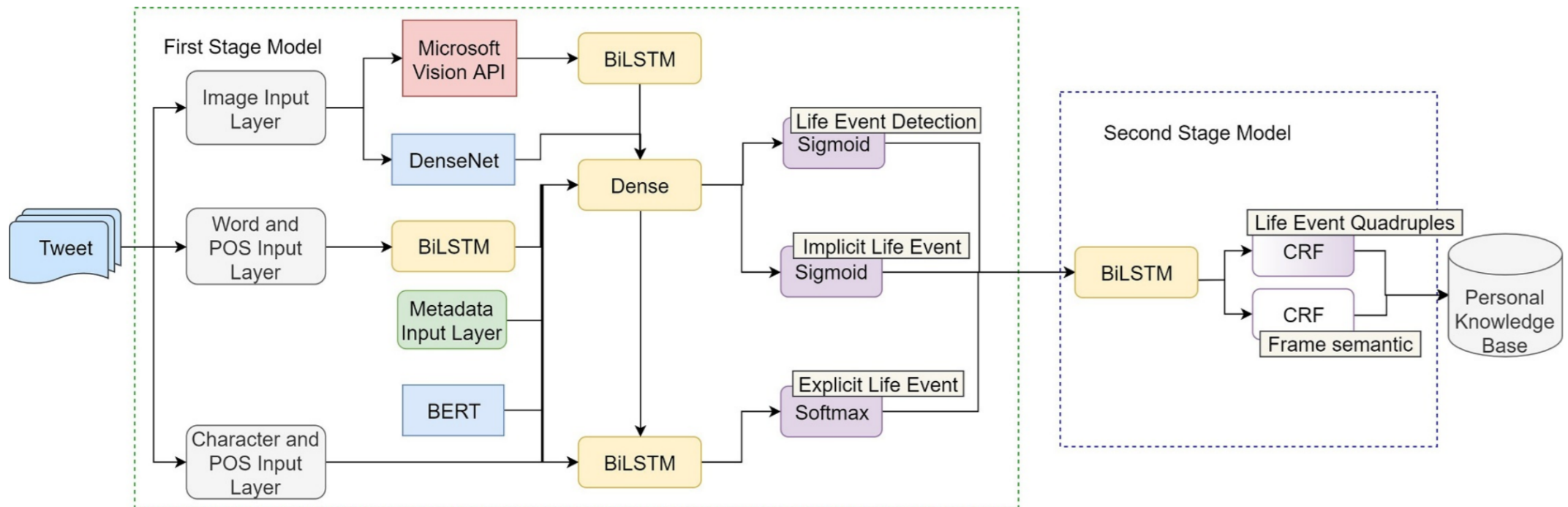
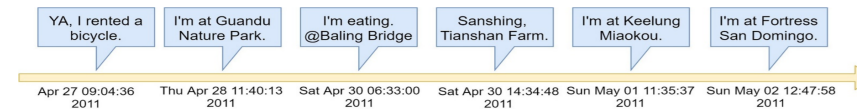
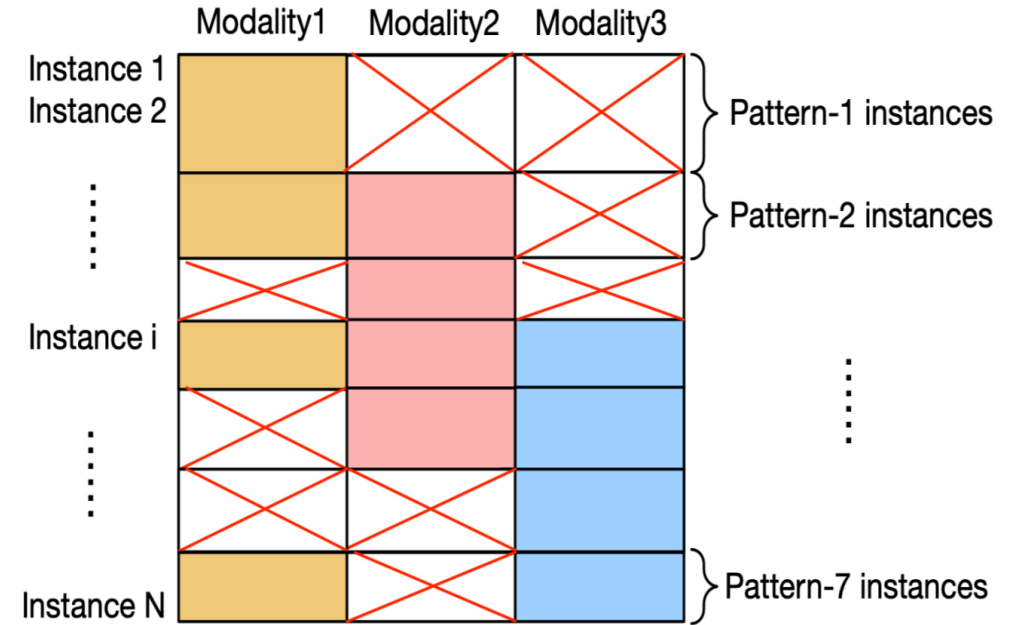
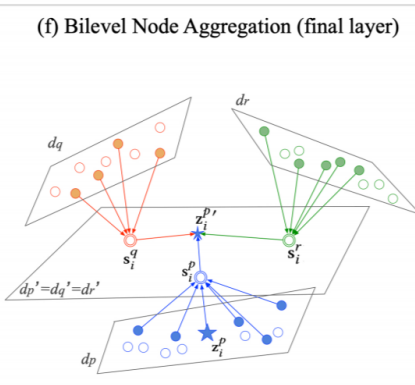
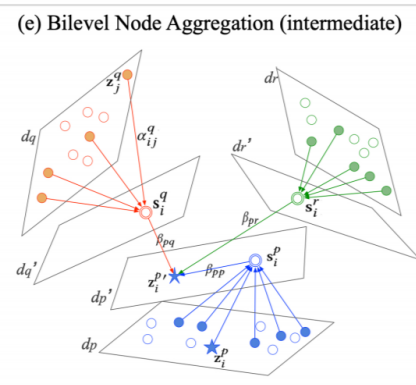
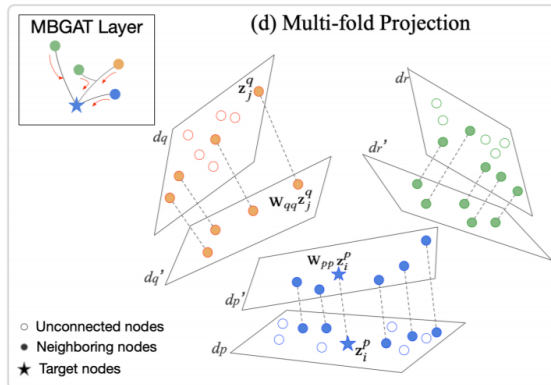
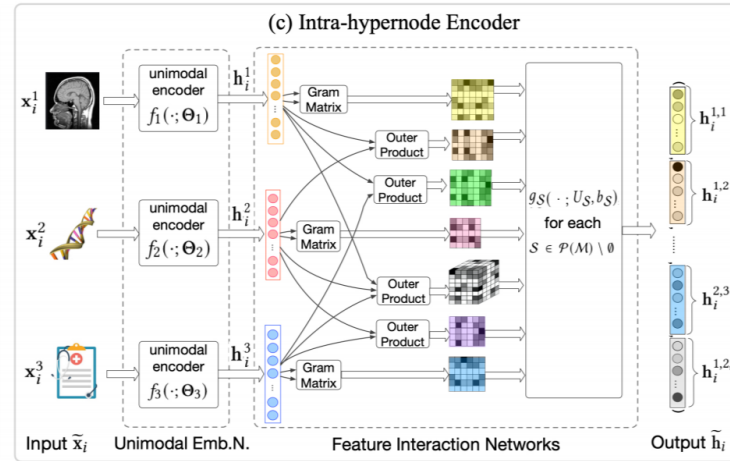
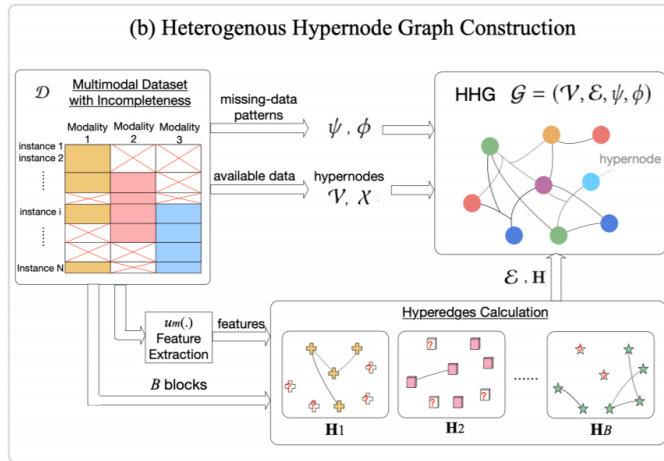
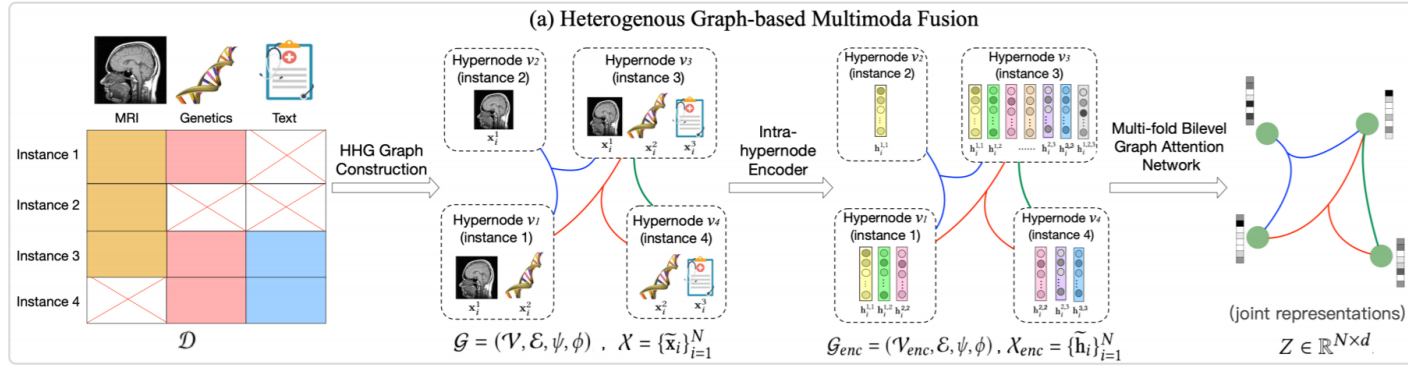


Fig. 1. An Example of social media post.



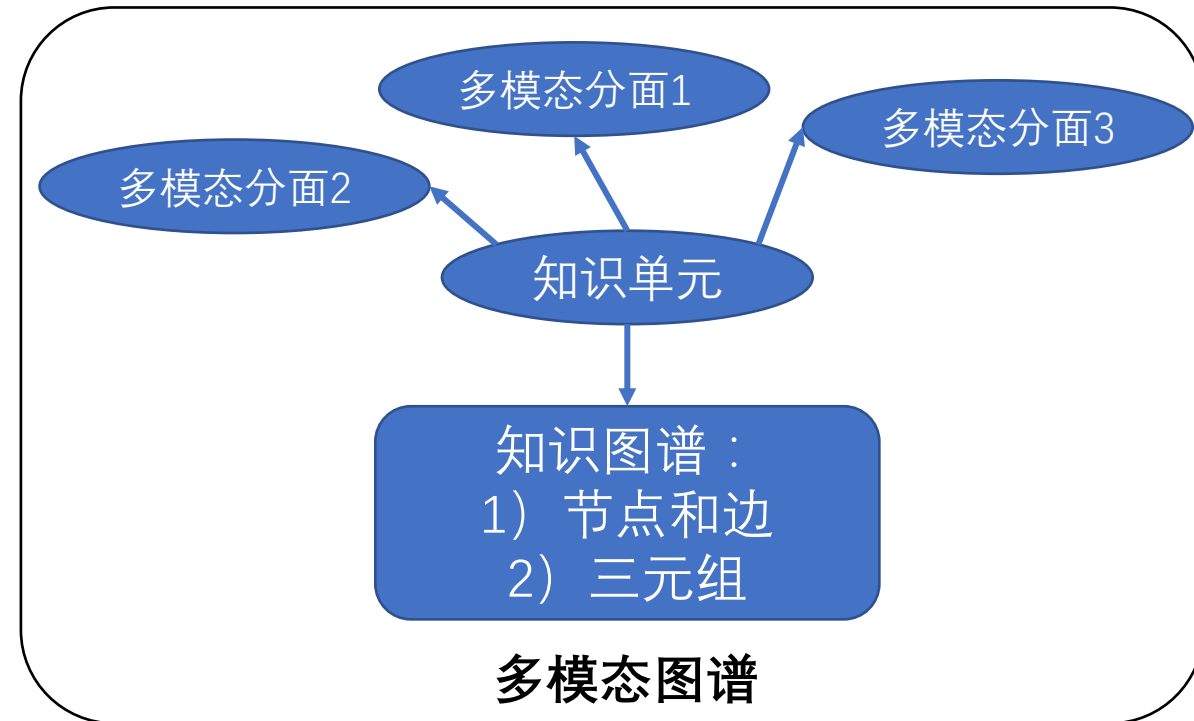
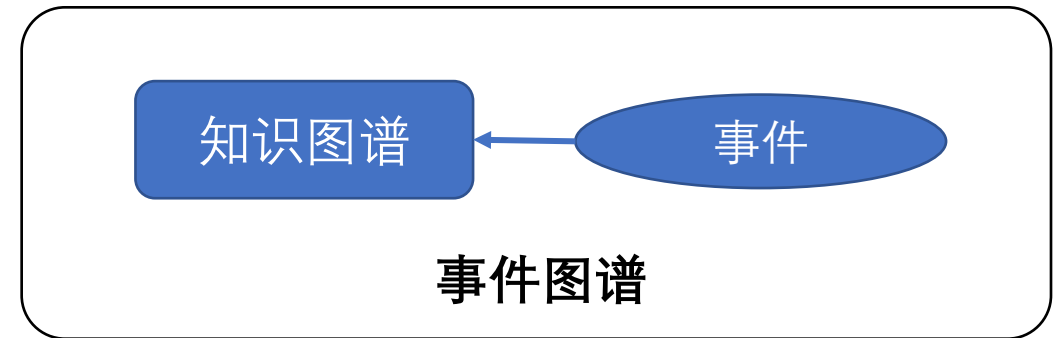
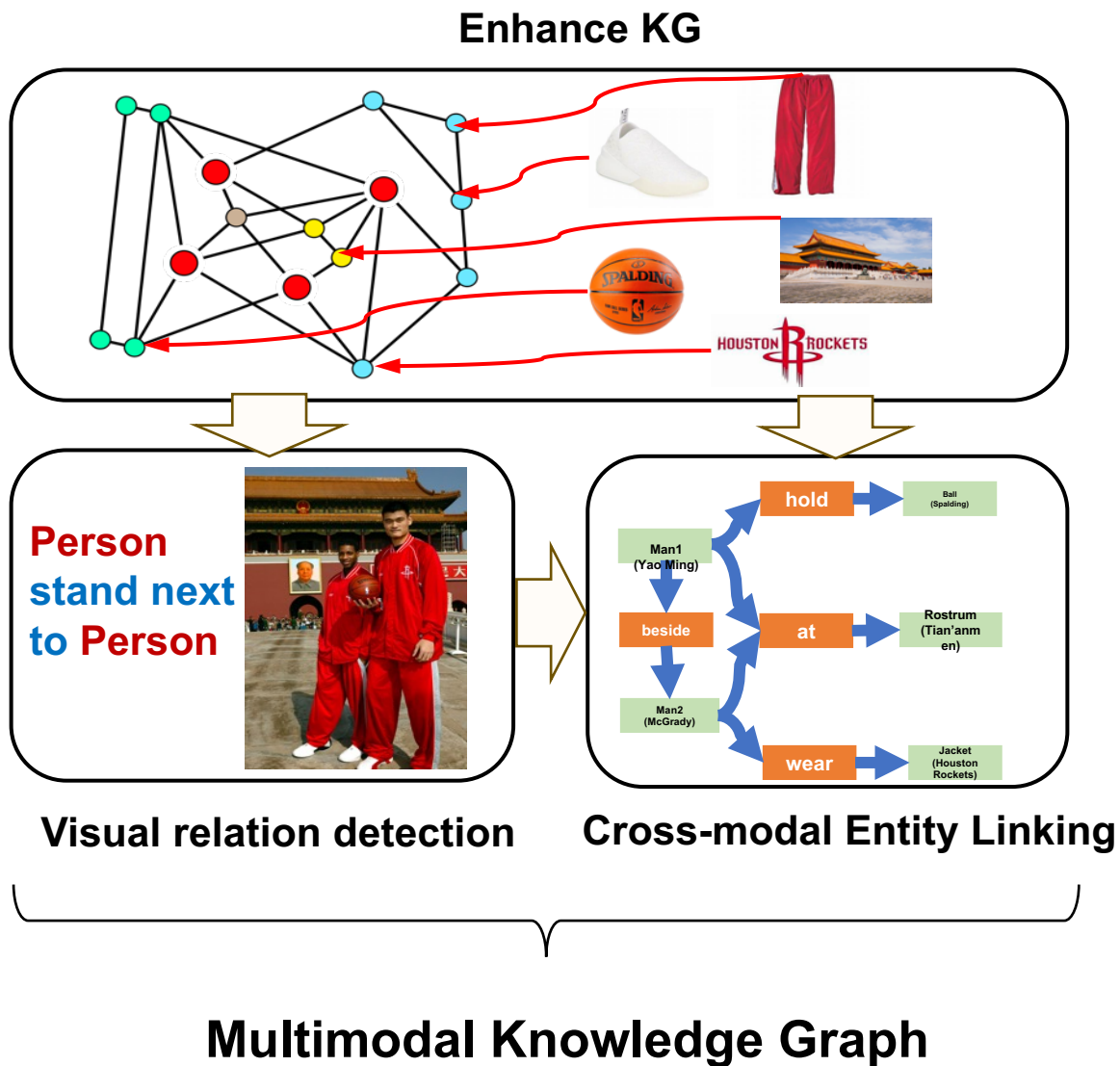
Yen, An-Zi, Hen-Hsen Huang, and Hsin-Hsi Chen. “Multimodal joint learning for personal knowledge base construction from Twitter-based lifelogs.” *Information Processing & Management* (2019): 102148.

Emerging Multimodal KG Work



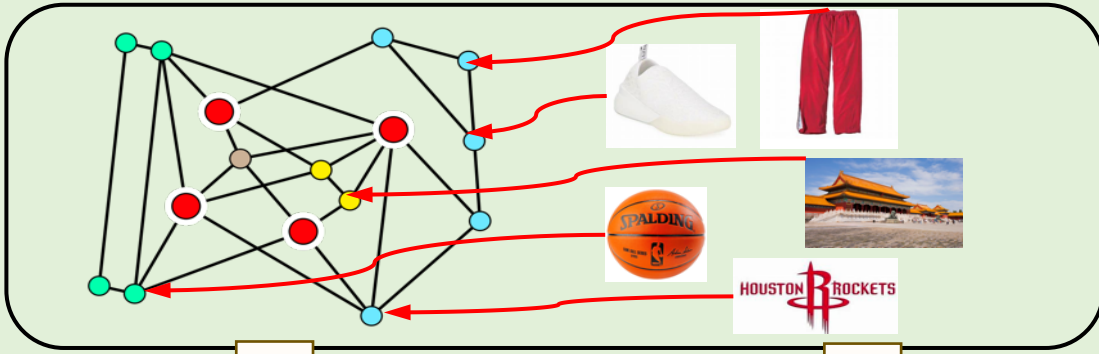
Jiayi Chen, Aidong Zhang “HGMF: Heterogeneous Graph-based Fusion for Multimodal Data with Incompleteness”, KDD 2020

Our Multimodal Knowledge Graph



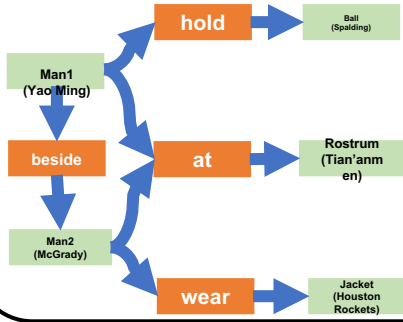
- Multimodality
- Multimodal KG Construction
- Inference
- Challenges

Enhance KG



Enhance KG

**Person
stand next
to Person**



Visual relation detection

Cross-modal Entity Linking

Multimodal Knowledge Graph

Application

1. Enhance KG

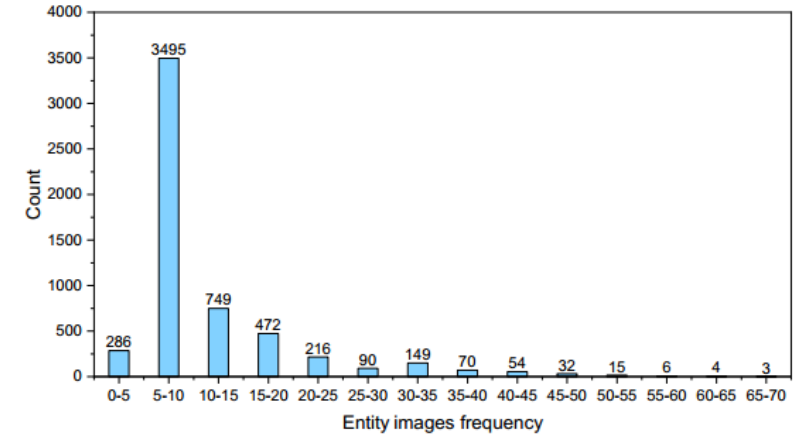
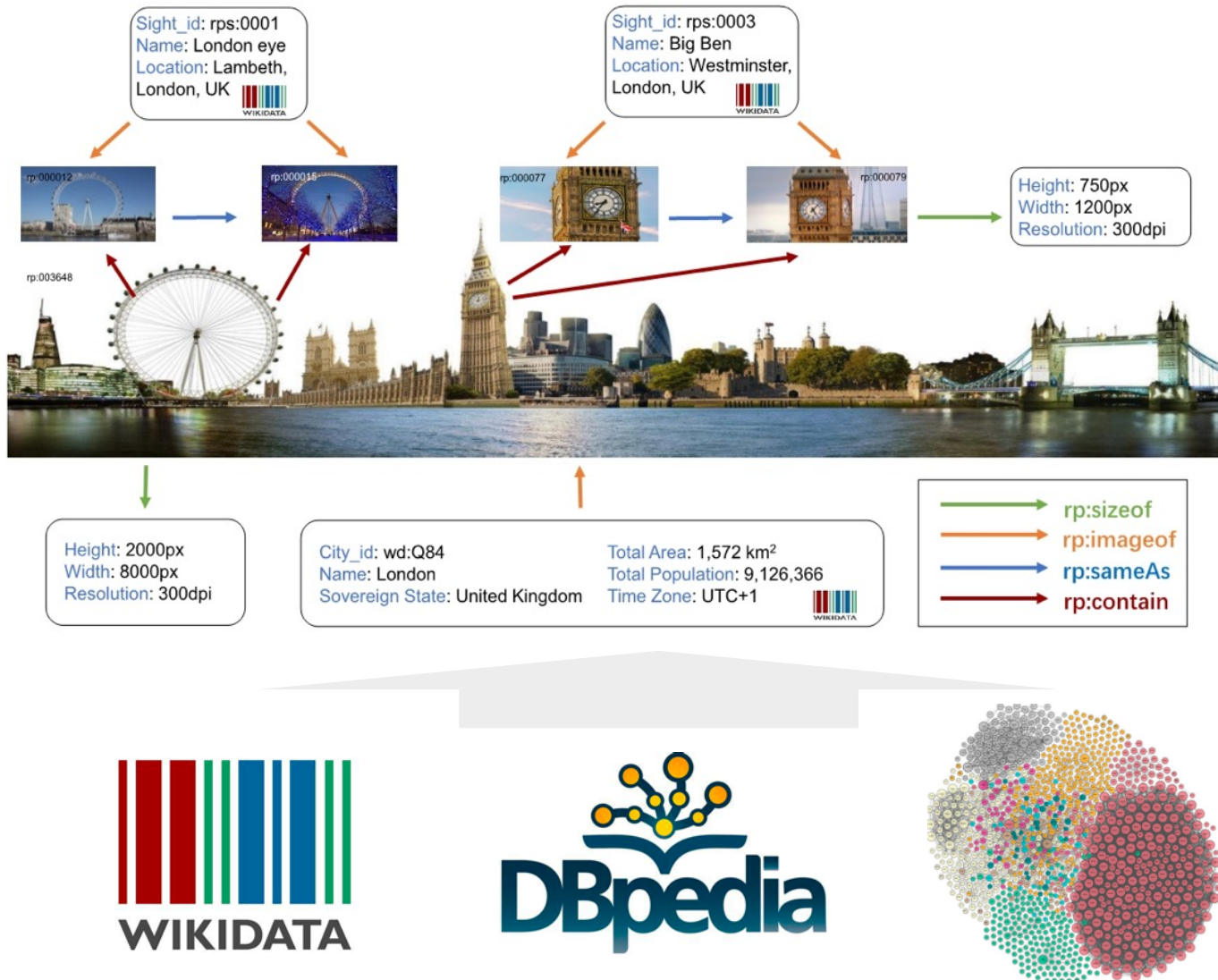
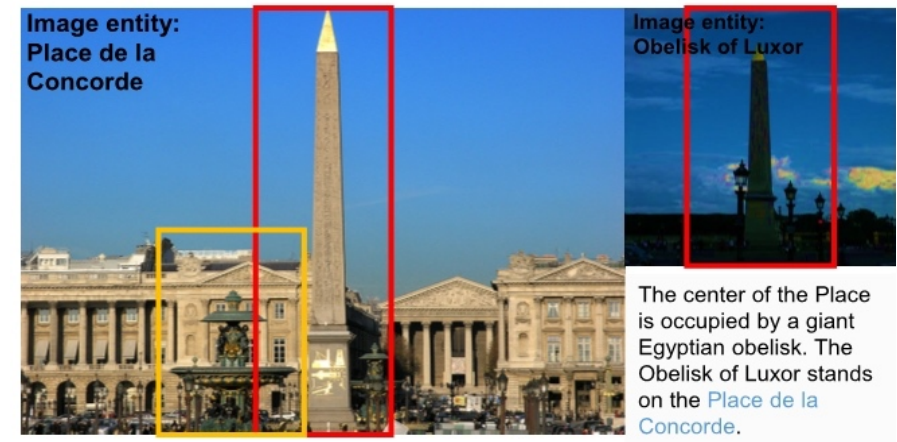
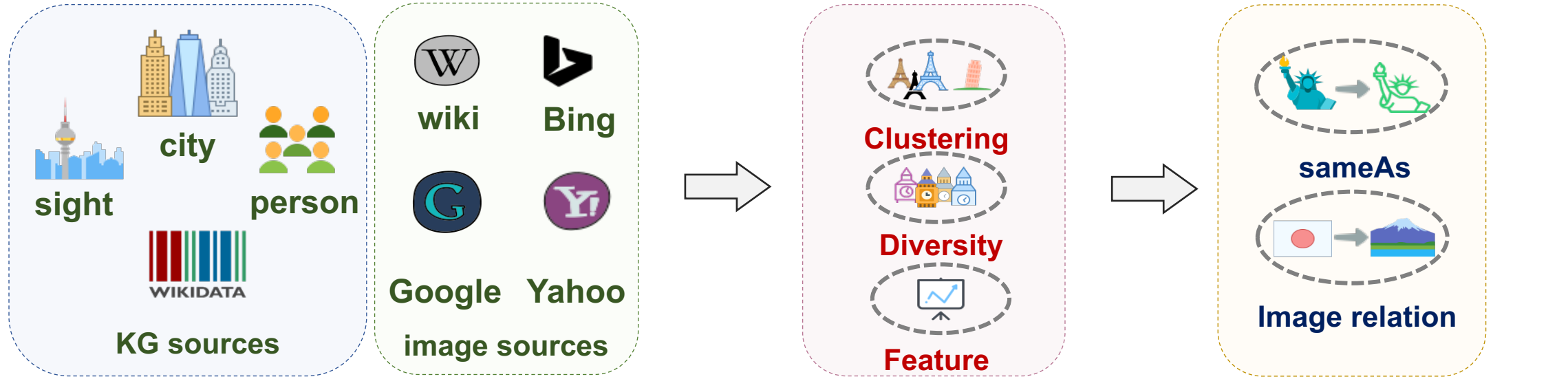


Image in KGs are very **sparse**



Missing visual relations

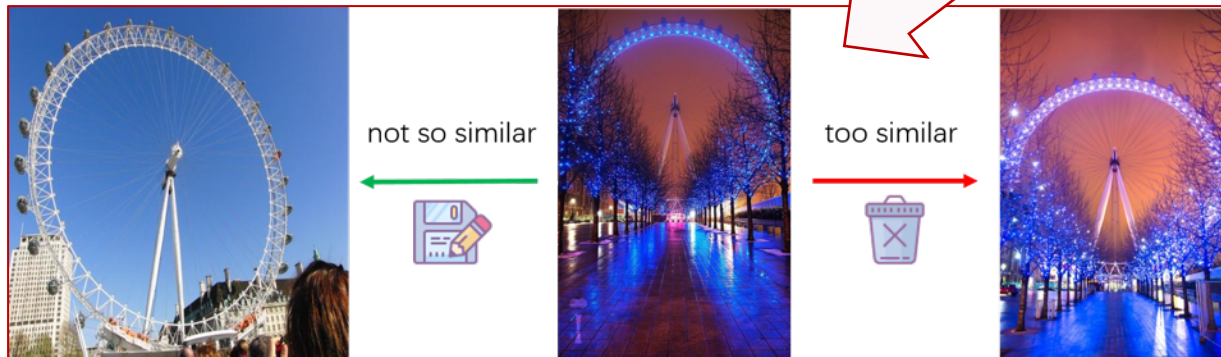
1. Enhance KG



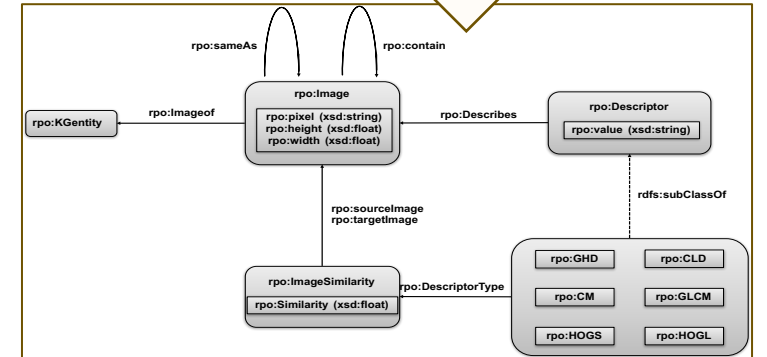
Data Collection

Image Processing

Relation Discovery



Diversity detection

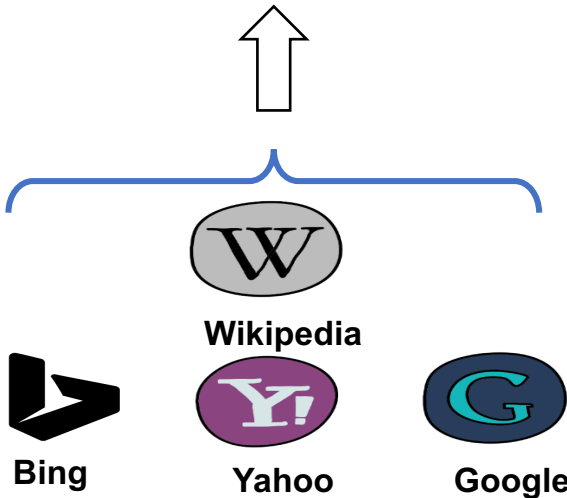
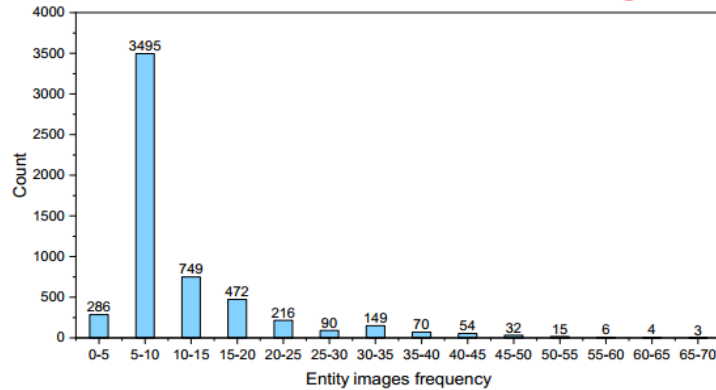


Visual relation ontology

1. Enhance KG

Richpedia is G , $G=(E, R)$, E : entities R : relations

Motivation: few images



Method: collect image resources

- Issue 1: Noise-containing image
- Method: K-means cluster
- Implement: Noise is an outlier in the clusters, and outliers are removed

- Issue 2: High image similarity
- Method: Diversity detection algorithm
- Implement:

$$\text{sim}(e_i, e_j) = \sum_{k=1}^n \min(H_k(e_i) - H_k(e_j))$$

First choose the center of mass, then choose the point farthest from the center of mass, and select the farthest point of those selected points in turn.

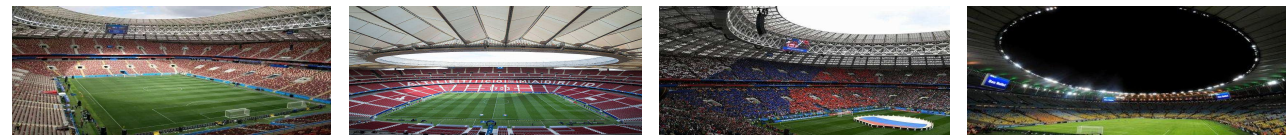
1. Enhance KG

Richpedia is G , $G=(E, R)$, E : entities R : relations



For each image, we generate five descriptors:

- **rpo:GHD**(Gradation Histogram Descriptor)
- **rpo:CLD**(Color Layout Descriptor)
- **rpo:CMD**(Color Moment Descriptor)
- **rpo:GLCM**(Gray-level co-occurrence matrix)
- **rpo:HOG**(Histogram of Oriented Gradient)

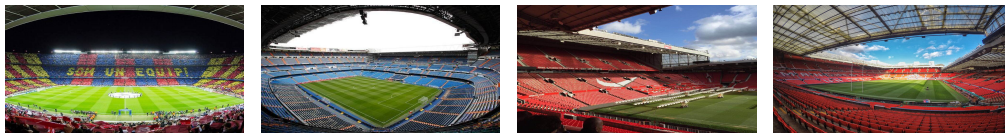


Source Image

1-NN

2-NN

3-NN



4-NN

5-NN

6-NN

7-NN



8-NN

9-NN

10-NN

1. Enhance KG

Richpedia is **G**, $G=(E, R)$, E: entities R: relations

City_id:wd:Q84
Name:London

rpo:imageof



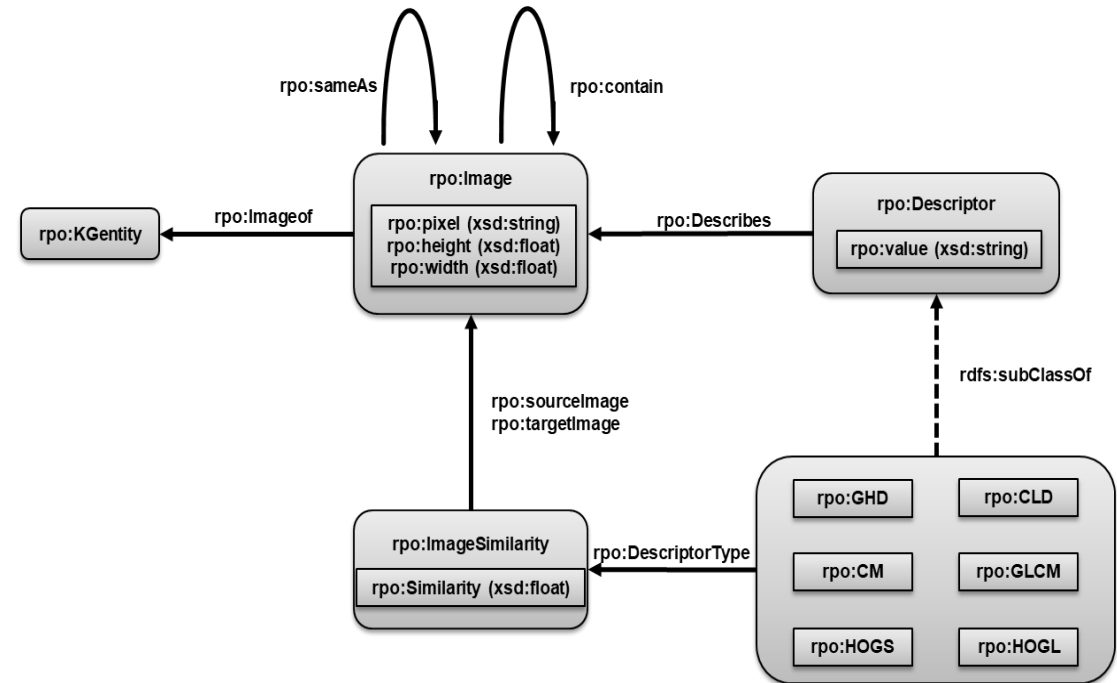
rpo:imageof



rpo:sameAs



<wd:Q84,wdt:P31,wd:Q515>
<rp:0000001,rpo:imageof,wd:Q84>
<rp:0000001,rpo:sameAs,rp:0000002>



Visual relation ontology

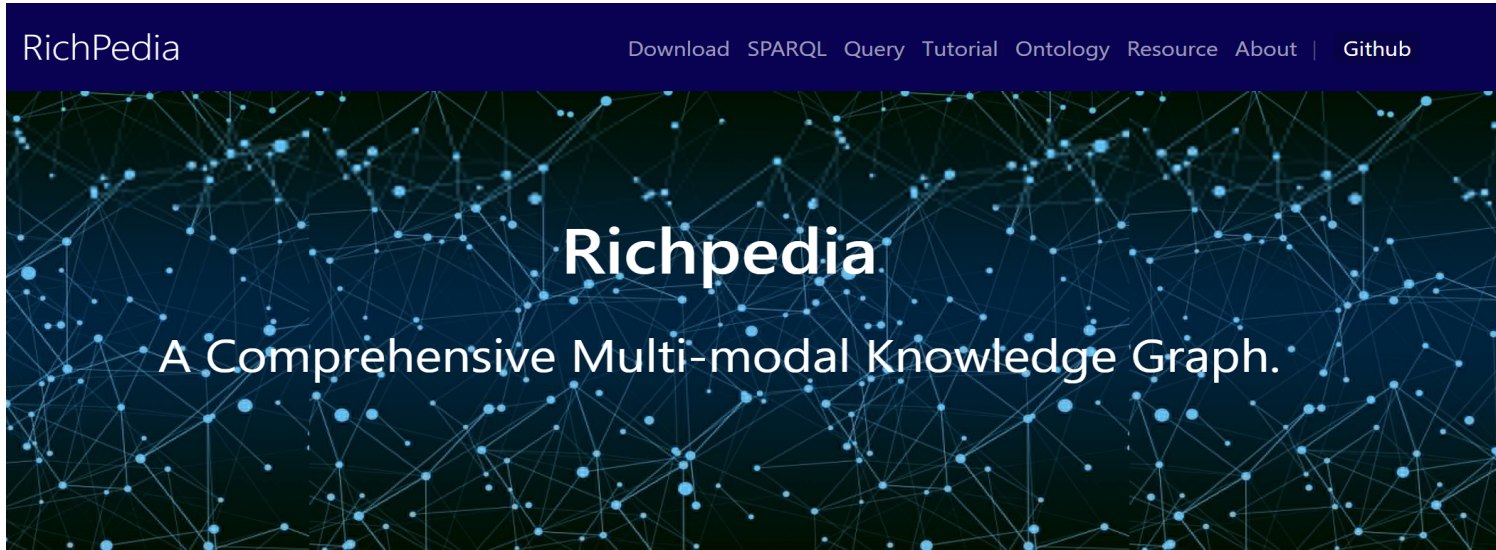
1. Enhance KG

Richpedia is G , $G=(E, R)$, E: entities R: relations



- Rule based: **simple but effective and efficient**
- Rule1: If there is **a hyperlink** in the description, the relationship is discovered by **a string mapping algorithm** between the keyword and the predefined relational ontology.
- Rule2: If there are **multiple hyperlinks** in the description, cyclically input the KG entity corresponding to the hyperlink, **simplifying** the situation to Rule1.
- Rule3: If there are **no hyperlinks** in the description, use the named entity recognition tool to find the KG entities and **simplify** the situation to Rule1 and Rule2.

1. Enhance KG



Richpedia.cn

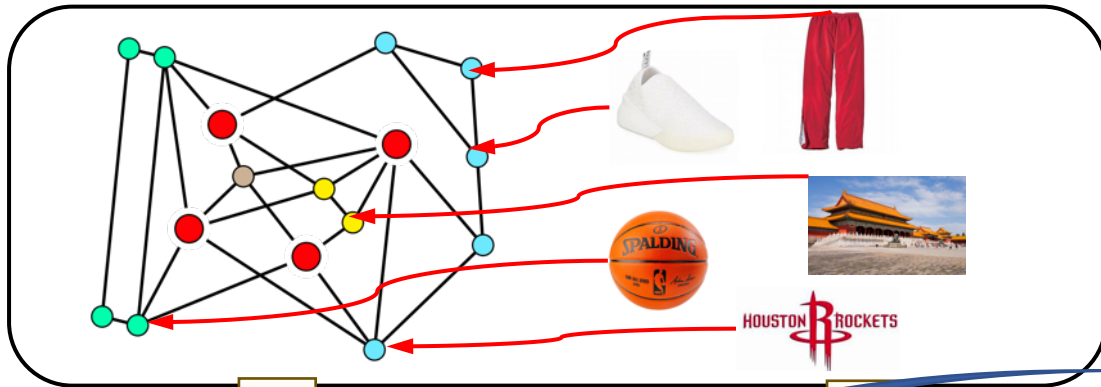
We will introduce how
to use this dataset later

Introduction

With the rapid development of Semantic Web technologies, various knowledge graphs are published on the Web using Resource Description Framework (RDF), such as Wikidata and DBpedia. Knowledge graphs provide for setting RDF links among different entities, thereby forming a large heterogeneous graph, supporting semantic search, question answering and other intelligent services. Meanwhile, public availability of visual resource collections has attracted much attention for different Computer Vision (CV) research purposes, including visual question answering, image classification, object and relationship detection, etc. And we have witnessed promising results by encoding entity and relation information of textual knowledge graphs for CV tasks. Whereas most knowledge graph construction work in the Semantic Web

**Meng Wang, Guilin Qi, Haofen Wang and
Qiushuo Zheng. Richpedia: A Large-Scale,
Comprehensive Multi-Modal Knowledge Graph.
Big Data Research, 2020**

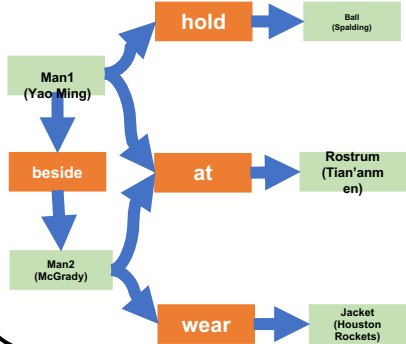
Enhance KG



Person
stand next
to
Person



Visual relation detection



Cross-modal Entity Linking

Visual Relation Detection

(Unbalanced, Long-tail)

Multimodal Knowledge Graph

Application

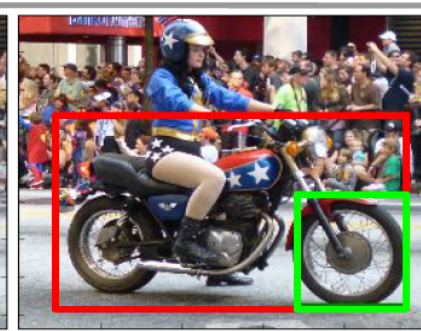
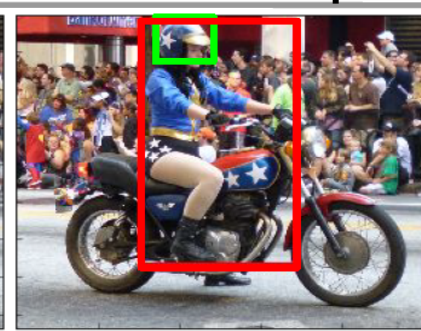
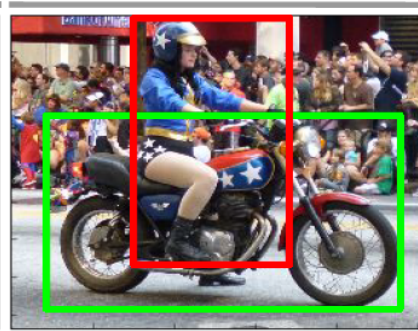
2. Visual Relation Detection

The aim of visual relation detection is to provide a comprehensive understanding of an image by describing all the objects within the scene, and how they relate to each other

Input



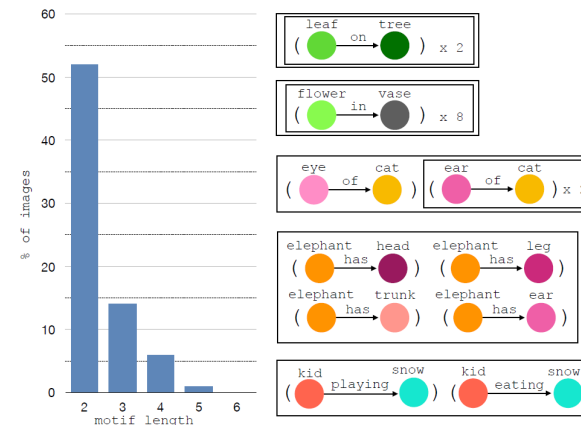
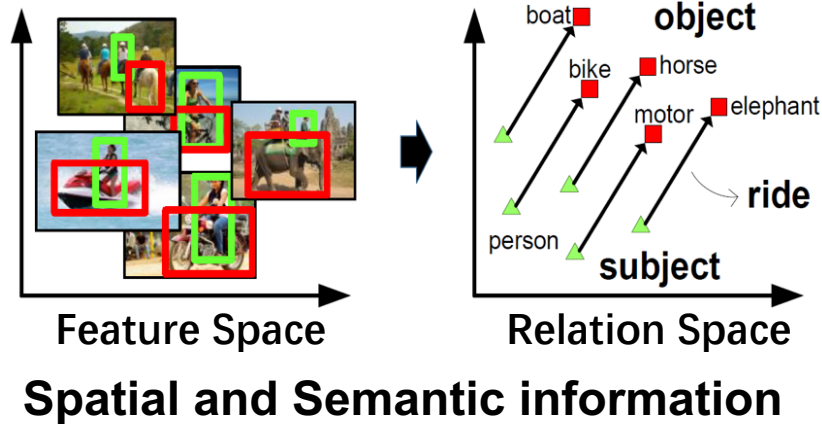
Output



person-on-motorcycle

person-wear-helmet

motorcycle-has-wheel



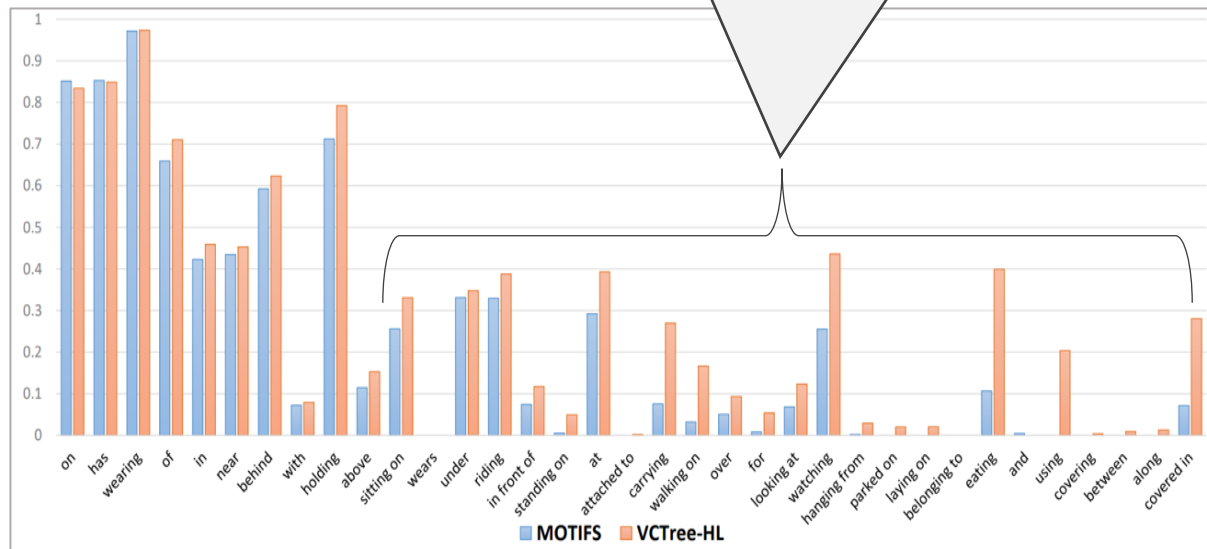
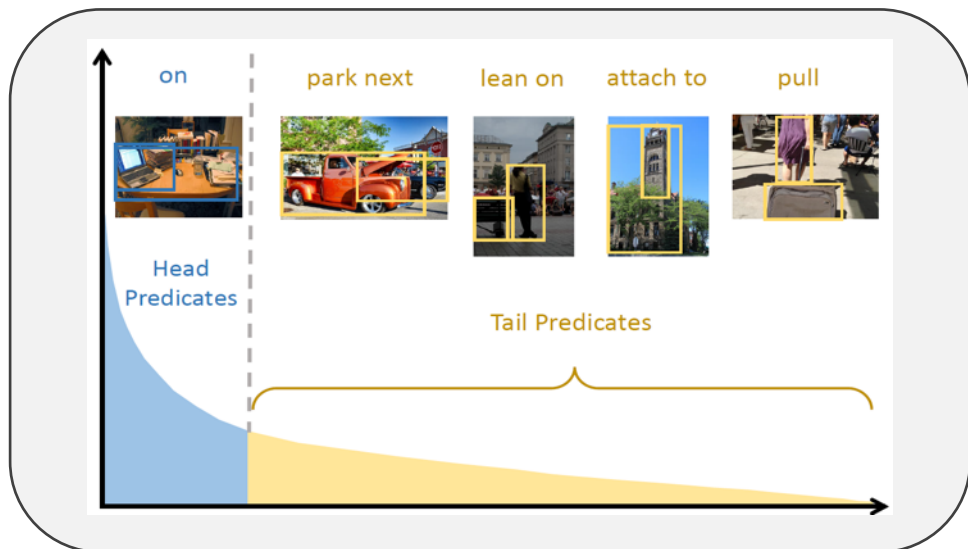
Statistical Information

VTransE, CVPR 2017

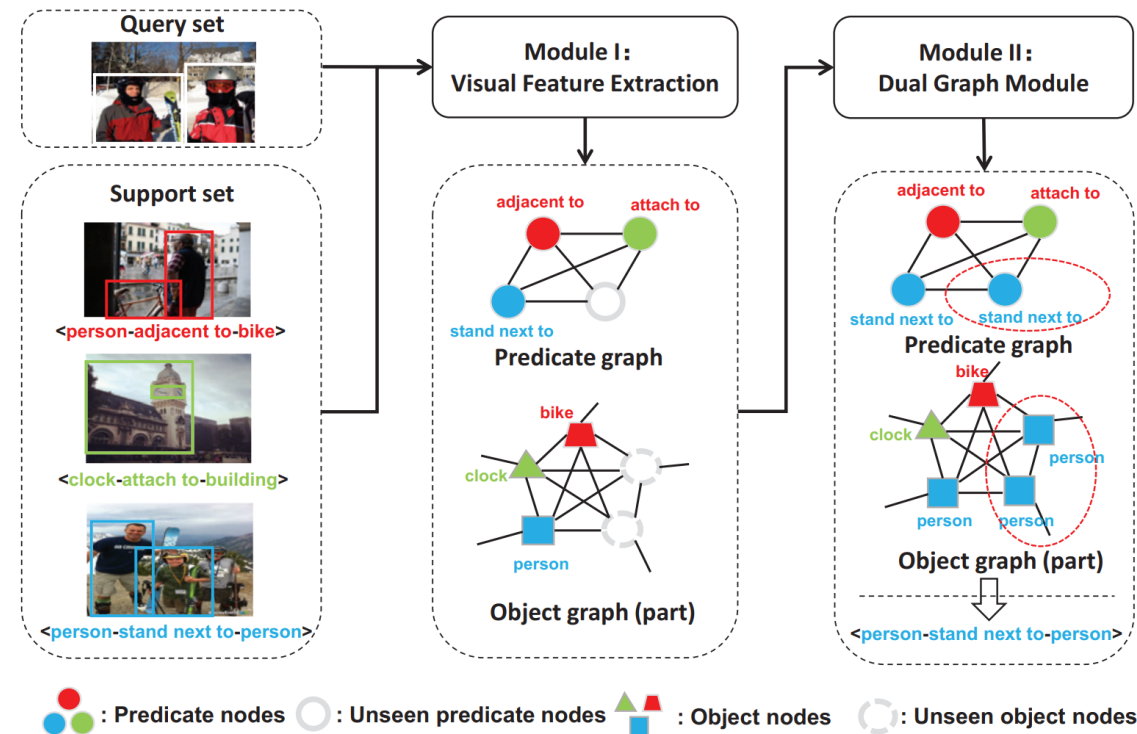
Language Priors, ECCV 2016

Neural motifs: Scene graph parsing with global context (CVPR2018)

2. Long-tail Visual Relation Detection

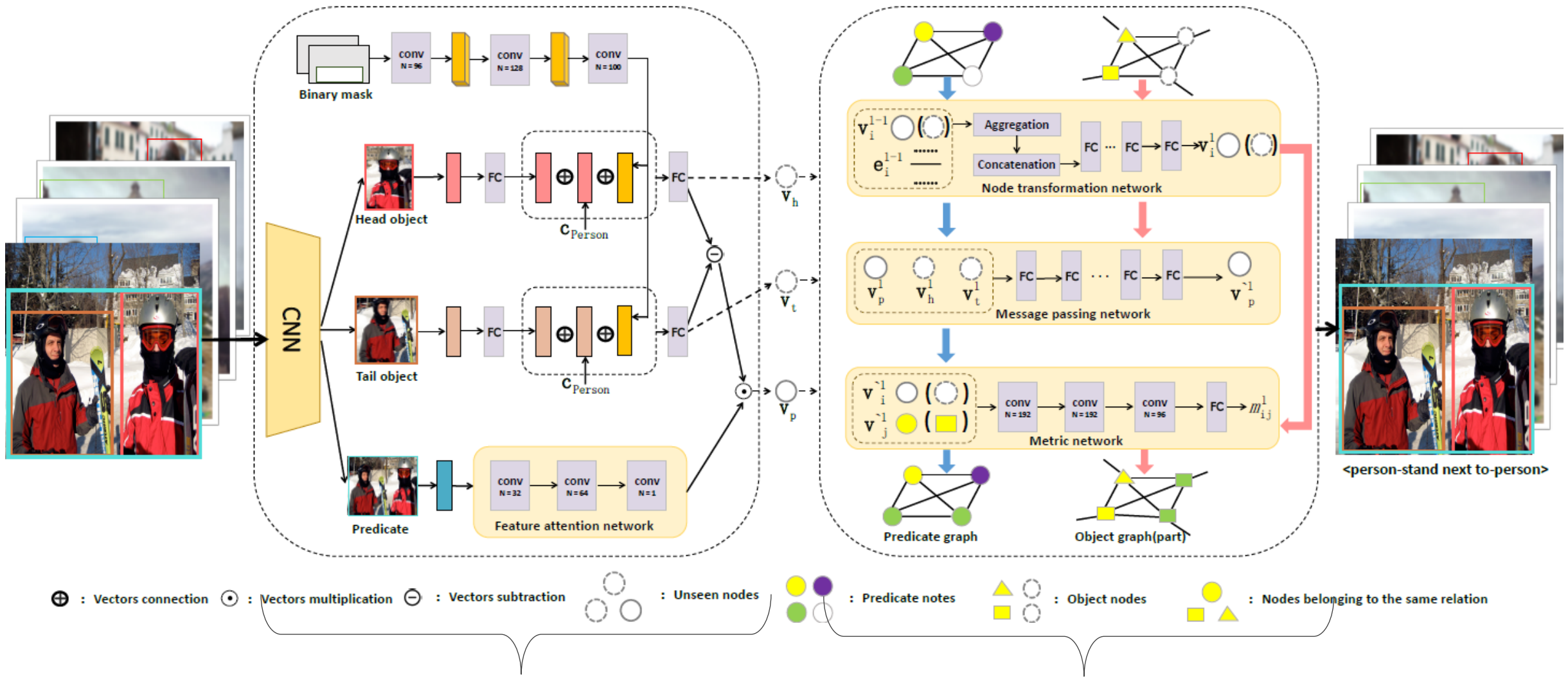


Learning to compose dynamic tree structures for visual contexts, CVPR 2019



Our method

2. Long-tail Visual Relation Detection



feature sparsity can be alleviated by using the feature-level attention mechanism

achieve messages passing from predicates or objects in different images

2. Long-tail Visual Relation Detection

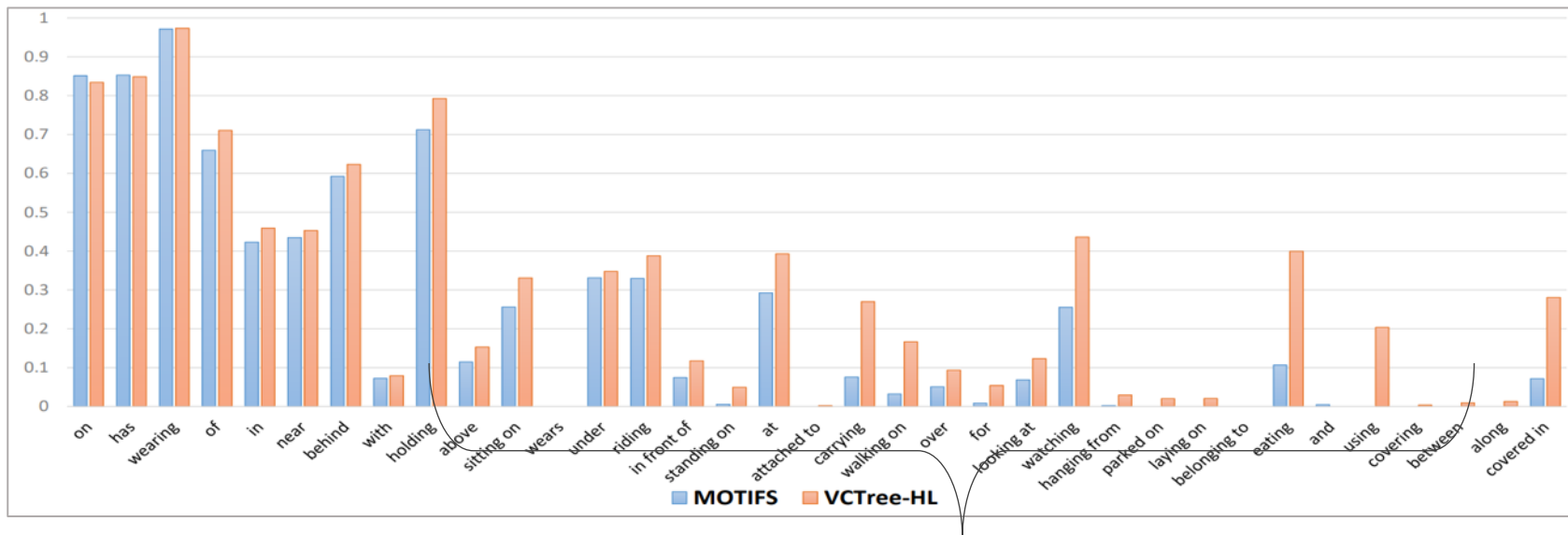


Table 2: Comparison with state-of-the-art baselines on the VRD-One and VG-One datasets.

	VRD-One				VG-One			
	PredCls		SGCls		PredCls		SGCls	
	5-way 1-shot	10-way 1-shot	5-way 1-shot	10-way 1-shot	5-way 1-shot	10-way 1-shot	5-way 1-shot	10-way 1-shot
VRD	37.4%	24.7%	14.7%	12.5%	41.4%	27.2%	10.8 %	9.6 %
VTransE	37.3%	24.3%	15.8%	13.4%	39.7%	23.4%	10.1%	9.4 %
LSVRU	40.3%	27.1%	16.9%	14.0%	43.4%	27.0%	10.7%	10.1%
RelDN	40.1%	26.4%	17.2%	14.3%	43.7%	28.3%	11.3%	10.1%
RelDN w/o sem	40.6%	27.3%	17.4%	14.9%	44.1%	28.2%	11.6%	10.4%
Ours	48.4%	33.5%	22.3%	20.9%	56.3%	37.5%	14.9%	13.2%

Table 4: Ablation studies on our model.

	5-way 1-shot	
	PredCls	SGCls
ours w/o Object Graph	47.3%	20.4%
ours w/o Message Passing Network	47.2%	21.8%
ours w/o Attention Network	45.7%	21.6%
ours All	48.4%	22.3%

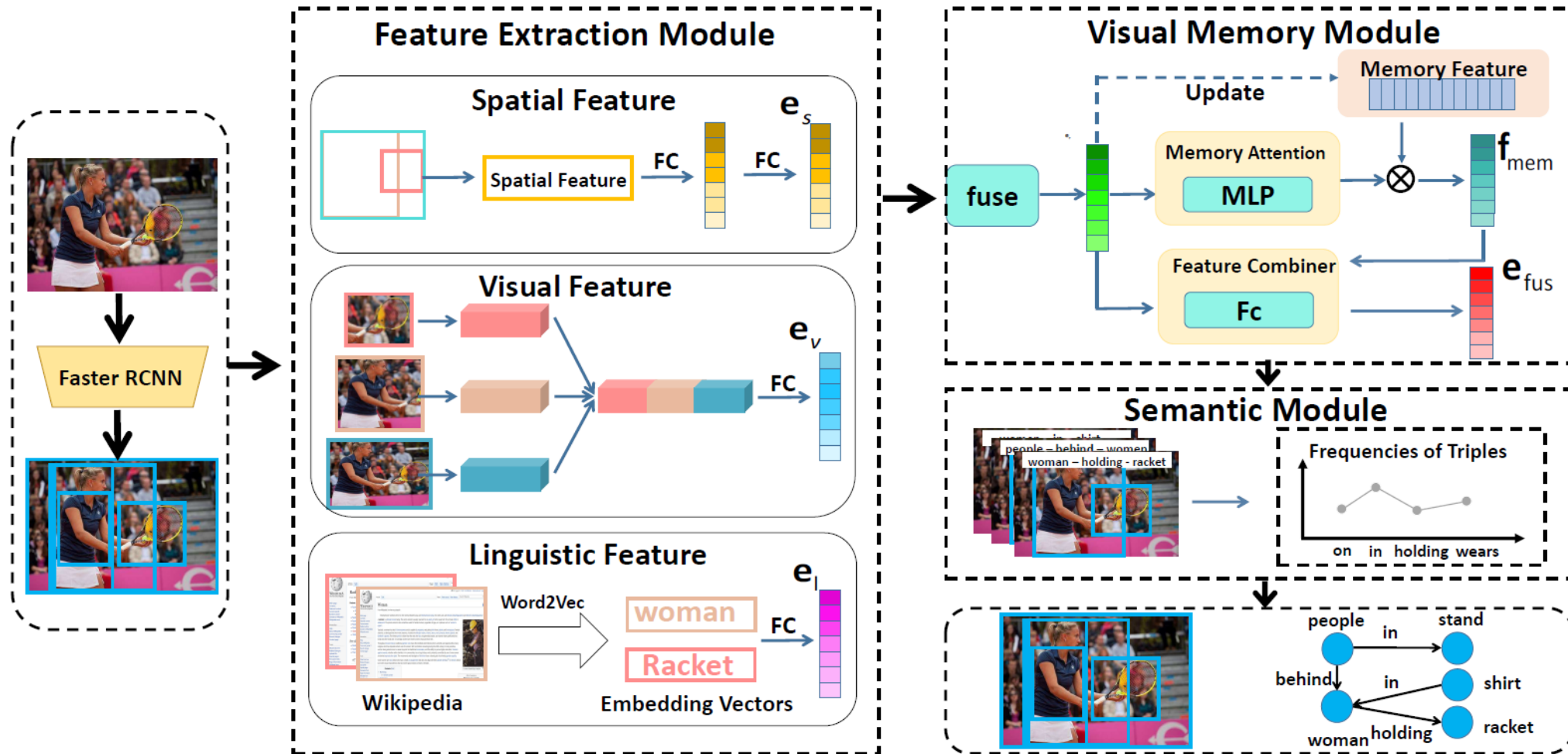
Weitao Wang, Meng Wang, Sen Wang, Guodong Long, Lina Yao, and Guilin Qi. One-Shot Learning for Long-Tail Visual Relation Detection. AAAI 2020.

2. Unbalanced Visual Relation Detection



Due to the existence of nonstandard labels, **excessive attention** to low-frequency visual relation will affect the performance of the scene graph generation model.

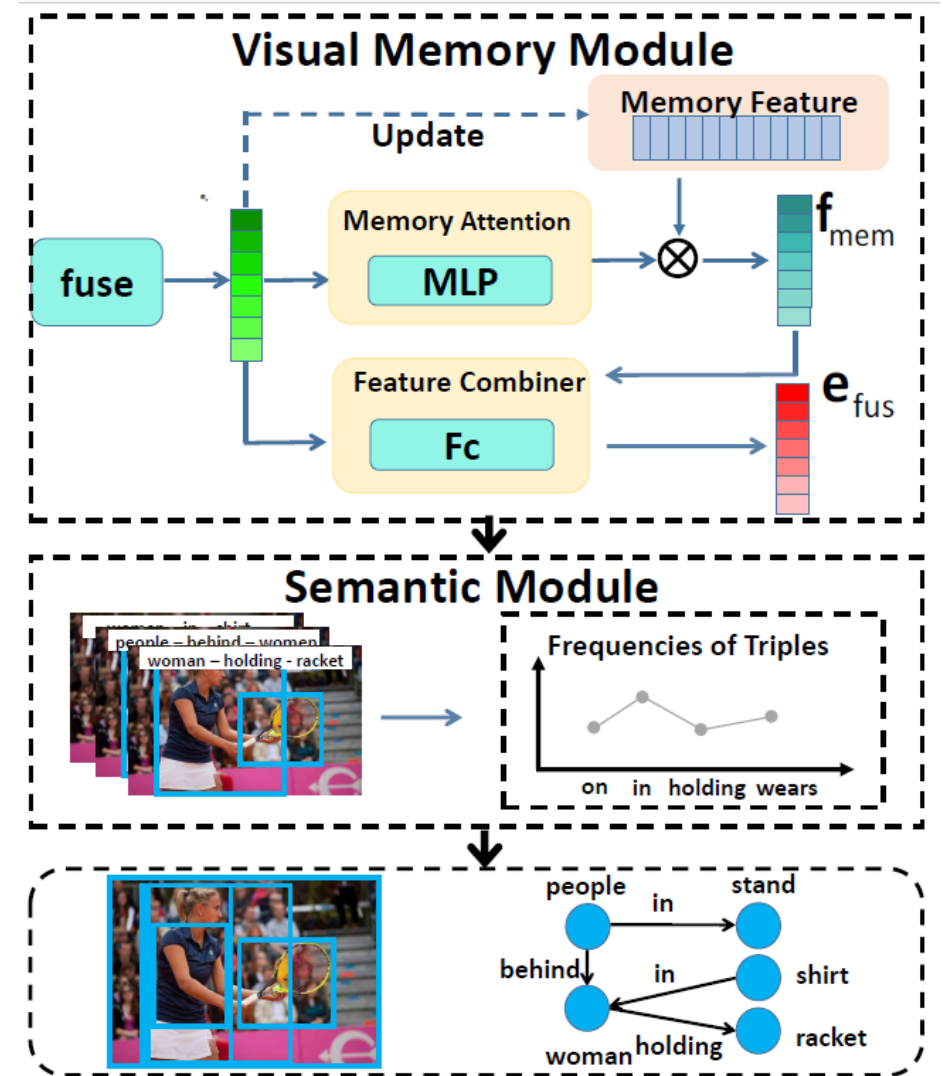
2. Unbalanced Visual Relation Detection



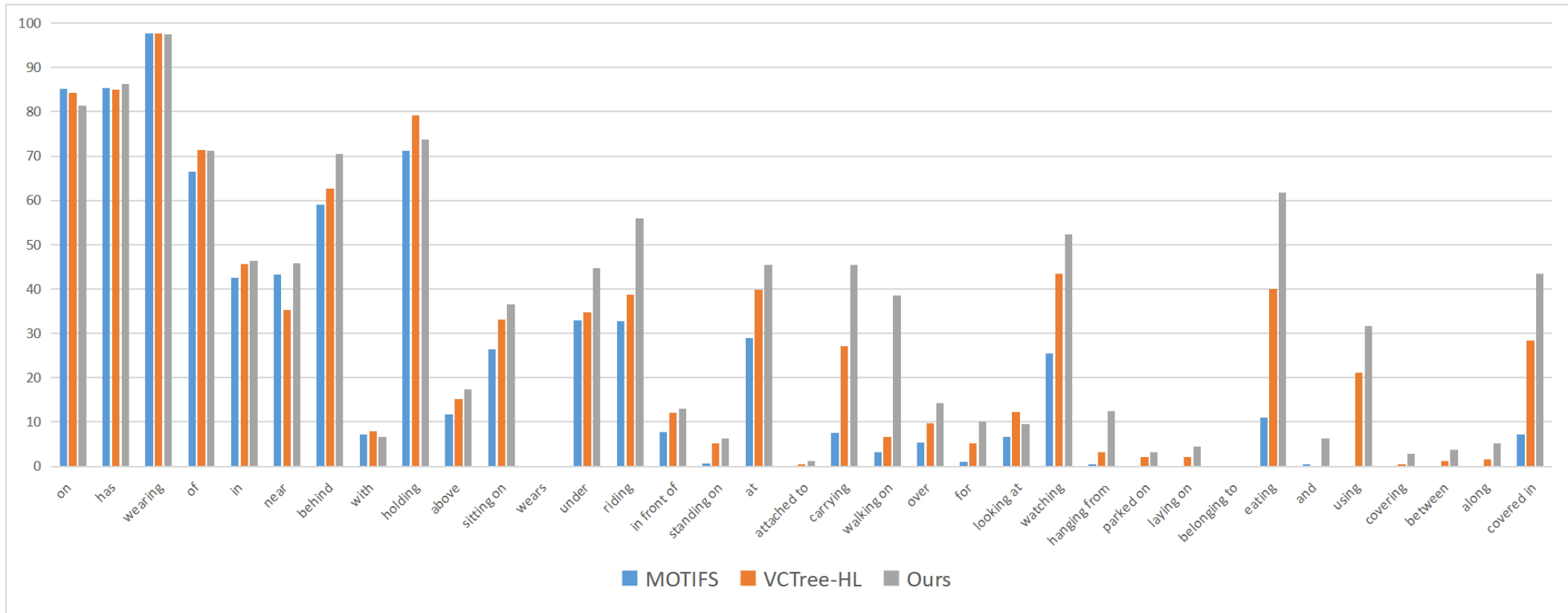
We use the method of memory features to realize the transfer of high-frequency relation features to low-frequency relation features.

2. Unbalanced Visual Relation Detection

- The calculation of visual relation memory is based on the prototype of each class, which is the mean of each category of features in the training set.
- The direct observation features and memory features are fused to realize the information exchange between the current relation and other relations
- We also utilize the statistical information (distribution) from the training set to influence the results of the model.



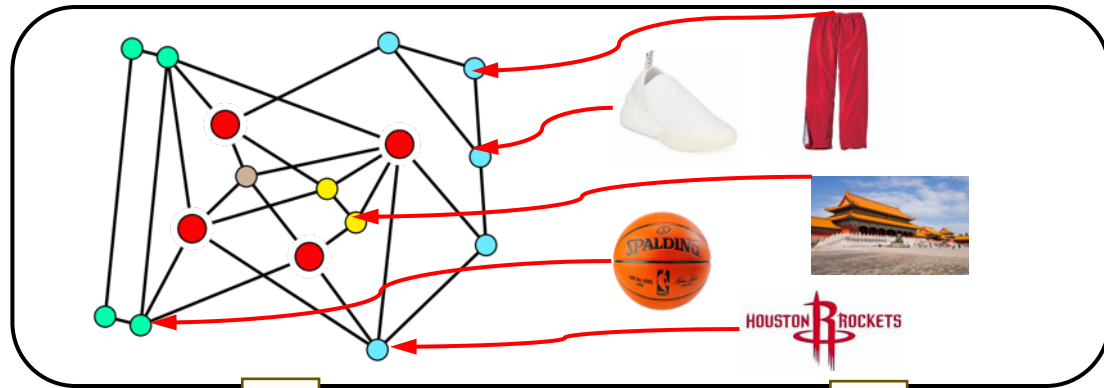
2. Unbalanced Visual Relation Detection



Our model achieves evident improvement in almost all relations (47/50)

Weitao Wang, Ruyang Liu, Meng Wang, Sen Wang, Xiaojun Chang, and Yang Chen.
Memory-Based Network for Scene Graph with Unbalanced Relations. 28th ACM MM 2020.

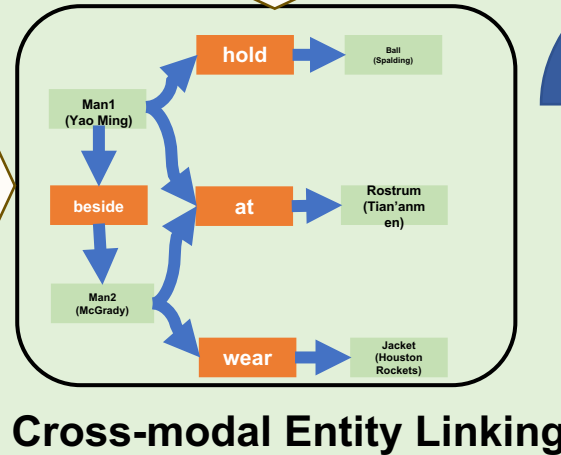
Enhance KG



Person
stand next
to
Person



Visual relation detection



Cross-modal Entity Linking

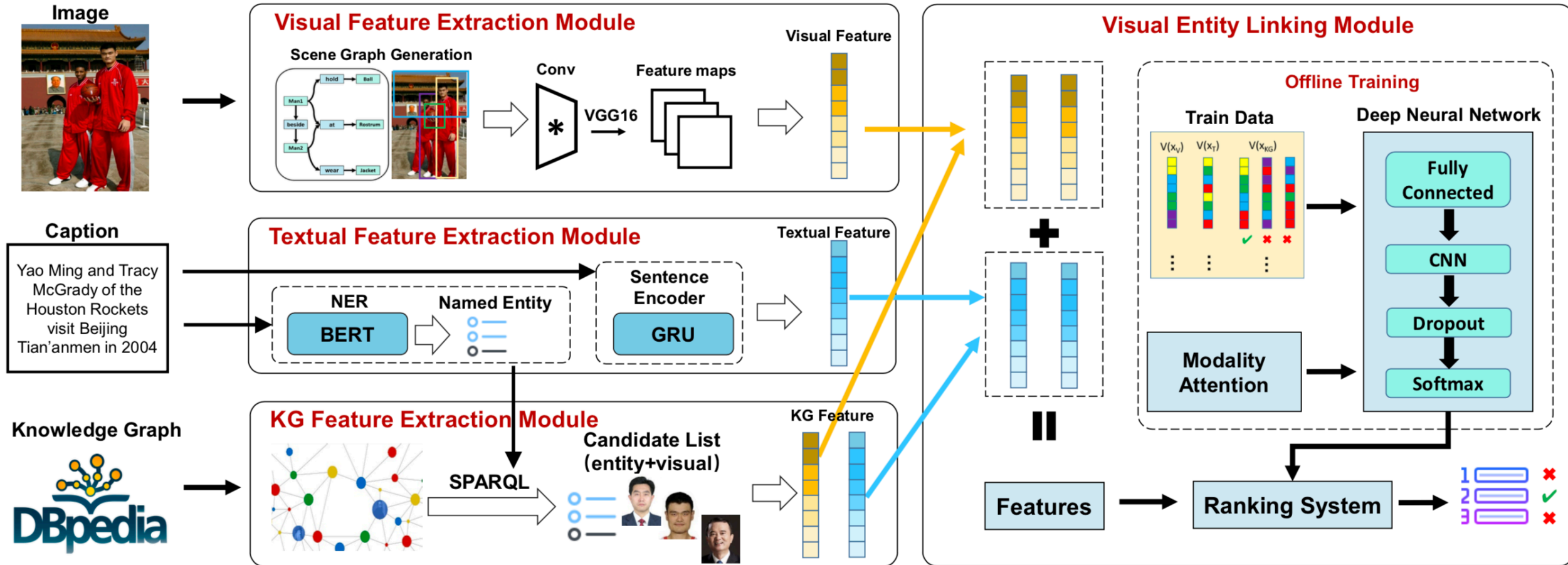
Multimodal Knowledge Graph

Application

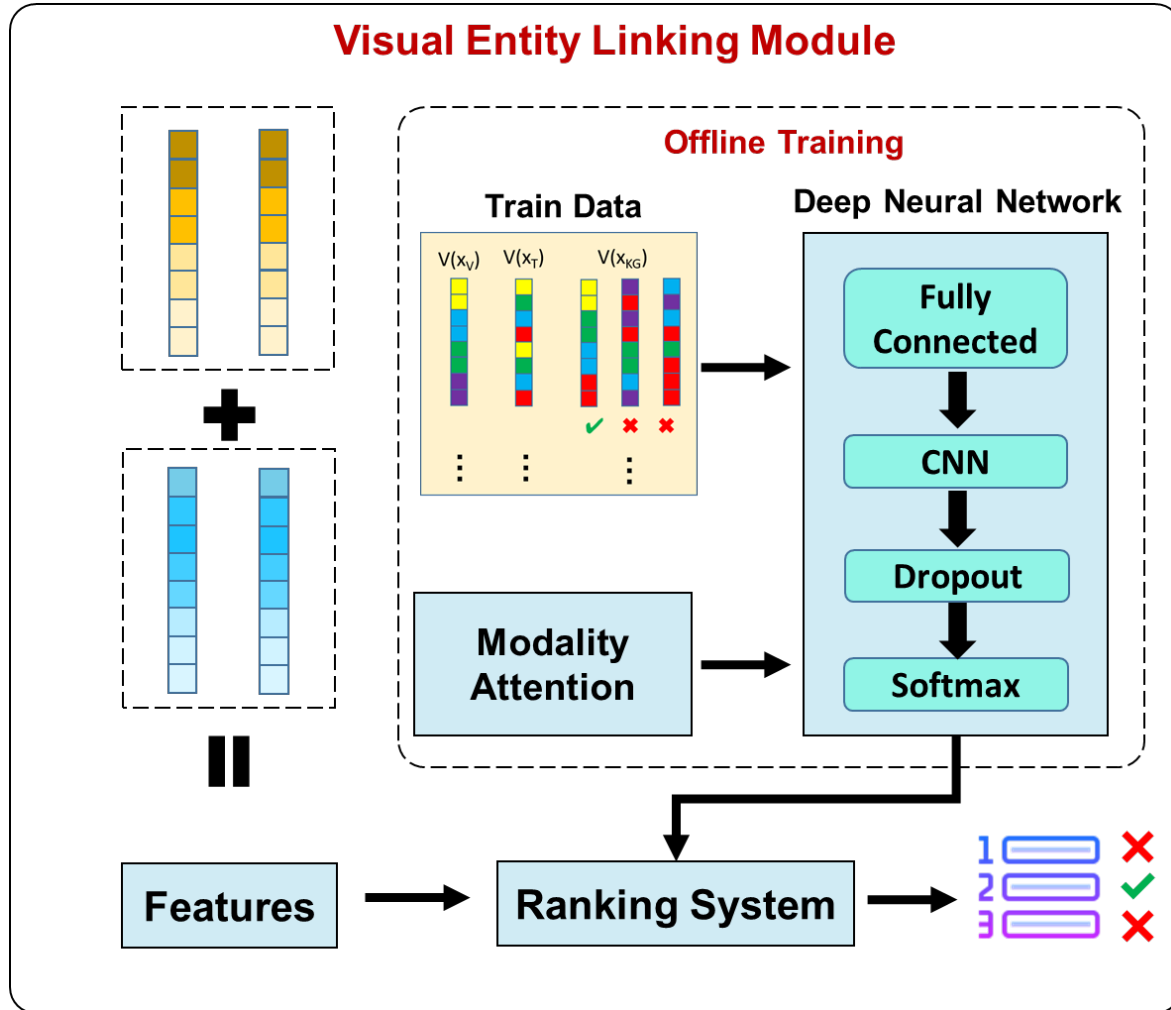
Cross-modal Entity Linking

(from classes to entities)

3. Cross-modal Entity Linking



3. Cross-modal Entity Linking



$$\mathcal{L} = \sum_{s=1}^m (\lambda_t \mathcal{L}_T(\bar{\mathbf{x}}^s) + \lambda_v \mathcal{L}_V(\bar{\mathbf{x}}^s)) + \|\mathbb{W}\|_2^2$$

$$\mathcal{L}_T(\bar{\mathbf{x}}) = \sum [\gamma + \text{conf}_t(\mathbf{y}') - \text{conf}_t(\mathbf{y})]_+$$

$$\mathcal{L}_V(\bar{\mathbf{x}}) = \sum [\gamma + \text{conf}_v(\mathbf{y}') - \text{conf}_v(\mathbf{y})]_+$$

$$\text{conf}_t(\mathbf{y}^i) = \frac{\exp(f(\mathbf{e}(\mathbf{x}_t), \mathbf{e}(\mathbf{y}_t^i)))}{\sum_{j \in \mathcal{C}} \exp(f(\mathbf{e}(\mathbf{x}_t), \mathbf{e}(\mathbf{y}_t^j)))}$$

$$\text{conf}_v(\mathbf{y}^i) = \text{cosine}(\mathbf{e}(\mathbf{x}_v), \mathbf{e}(\mathbf{y}_v^i))$$

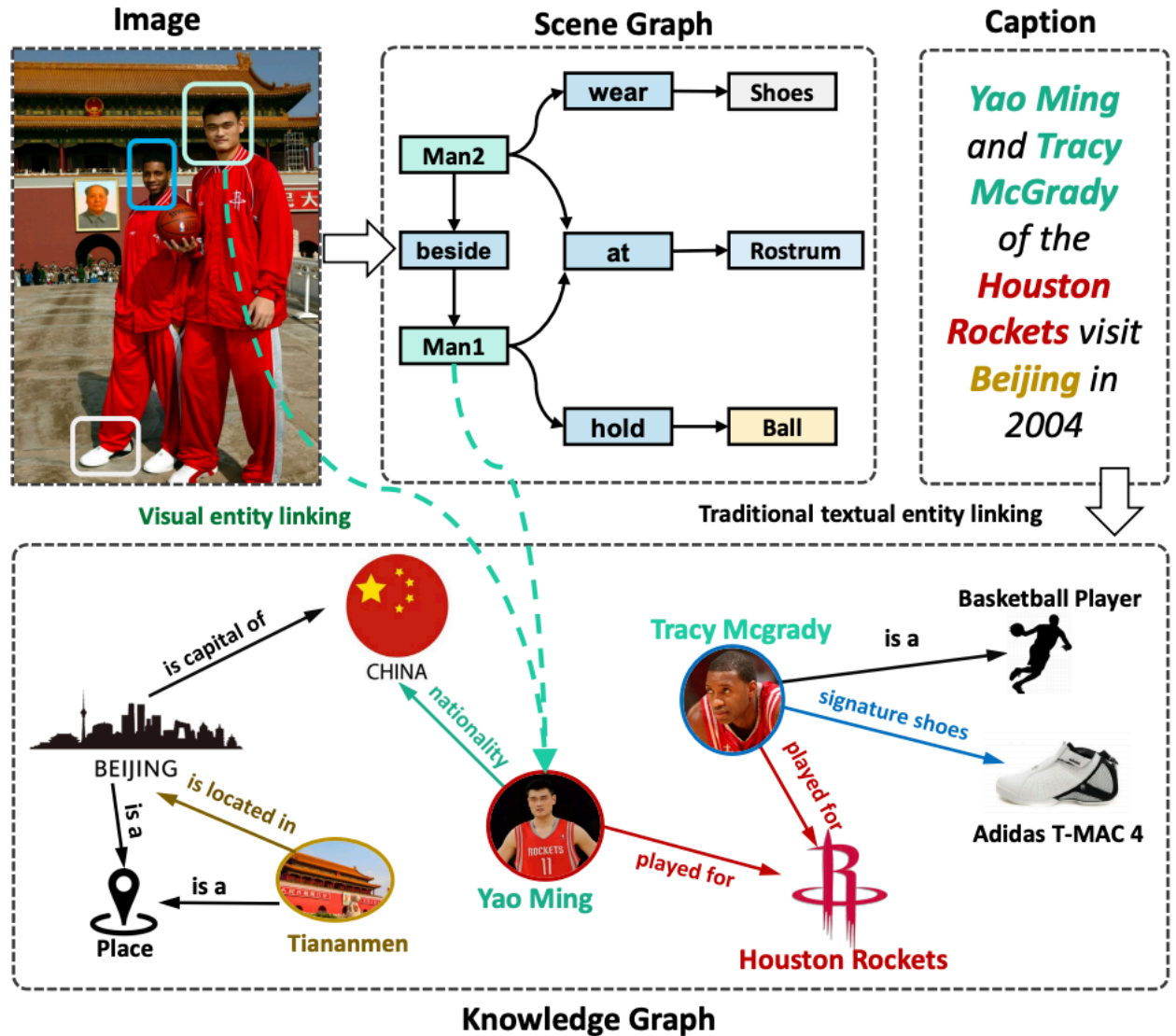
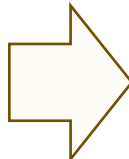
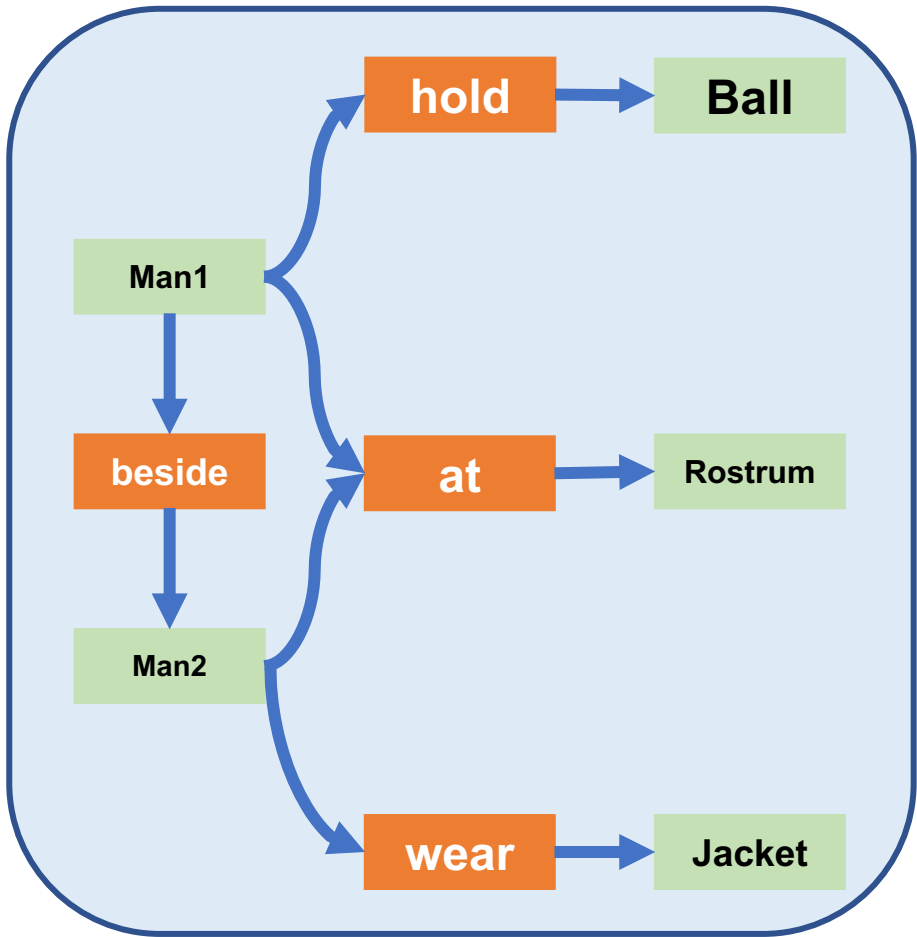
Learning to Rank With Modal Attention

$$[\mathbf{a}_t; \mathbf{a}_v] = \sigma(\mathbf{W} \cdot [\mathbf{x}_t; \mathbf{x}_v] + \mathbf{b})$$

$$\alpha_m = \frac{\exp(\mathbf{a}_m)}{\sum_{m' \in \{t, v\}} \exp(\mathbf{a}_{m'})} \quad \forall m \in \{t, v\}$$

$$\bar{\mathbf{x}} = \sum_{m \in \{t, v\}} \alpha_m \mathbf{x}_m$$

3. Cross-modal Entity Linking



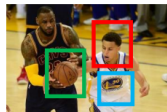
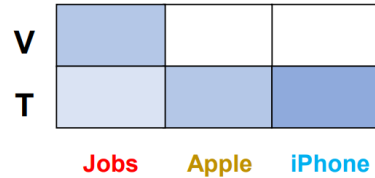
3. Cross-modal Entity Linking

	Accuracy	
	Top-1	Top-10
ours w/o visual features	56.19%	63.42%
ours w/o textual features	72.55%	82.23%
ours with a smaller sized KG	60.19%	66.47%
ours All	83.16%	93.81%

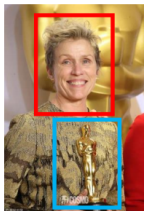
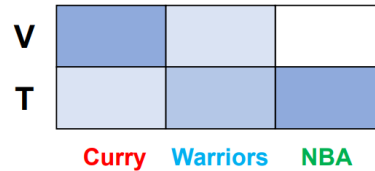
Modalities	Model	Accuracy			
		Top-1	Top-3	Top-5	Top-10
V+KG	Huawei API	12.53%	18.49%	20.46%	22.94%
V+KG	Tencent API	11.79%	16.42%	21.64%	24.61%
T+KG	Faster-RCNN+CoAtt	55.45%	63.76%	66.05%	67.91%
T+KG	Faster-RCNN+Falcon	56.16%	61.47%	62.17%	63.94%
T+KG	Faster-RCNN+CDTE	58.27%	64.79%	65.09%	66.14%
V+T+KG	DZMNED	66.46%	73.16%	81.06%	83.49%
V+T+KG	Our method	83.16%	88.61%	92.49%	93.81%



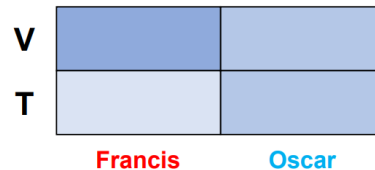
Jobs, **Apple's** founder, attended the launch of the new **iPhone**.



Curry won the **NBA** Championship for the **Golden State Warriors** at Auckland Stadium.



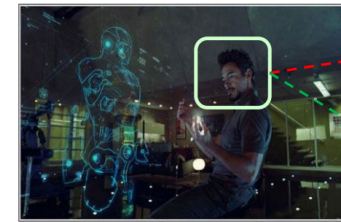
Francis McDonald win the **Best Actress Oscar** in 2018 and attended the awards ceremony.



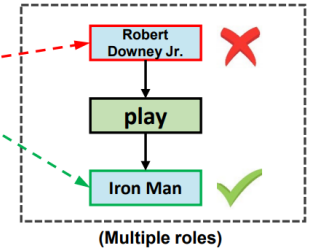
Caption

Robert Downey Jr. plays **Iron Man** in the movie

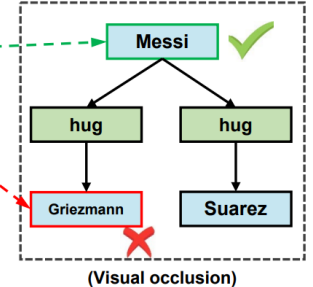
Image



Results and Ground Truth

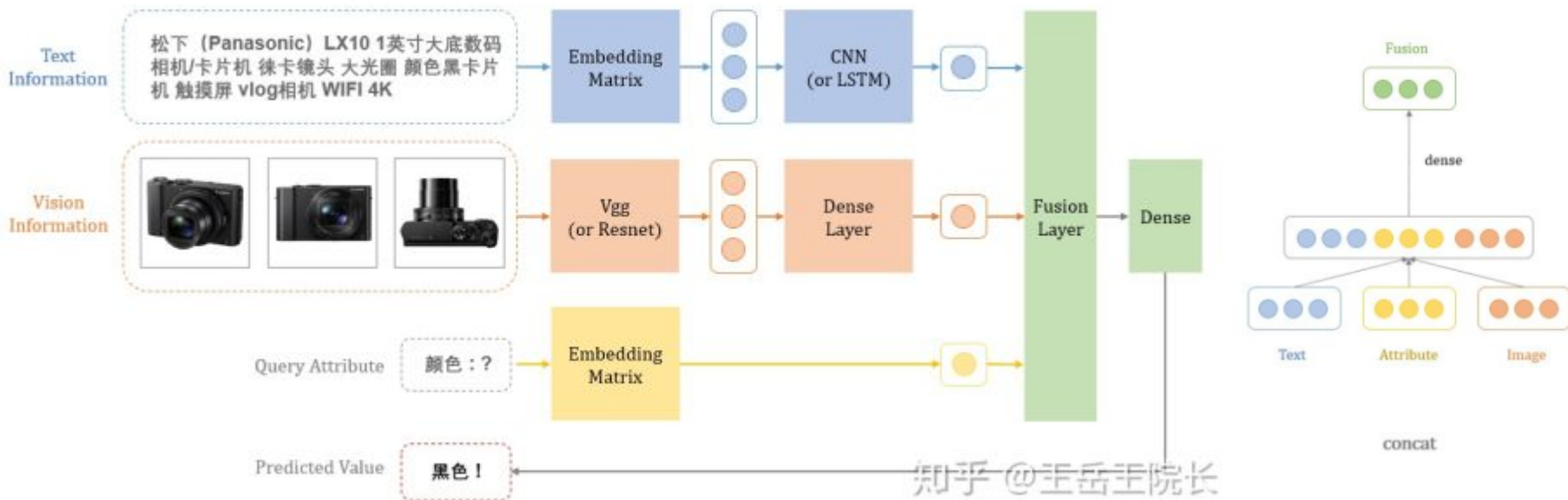
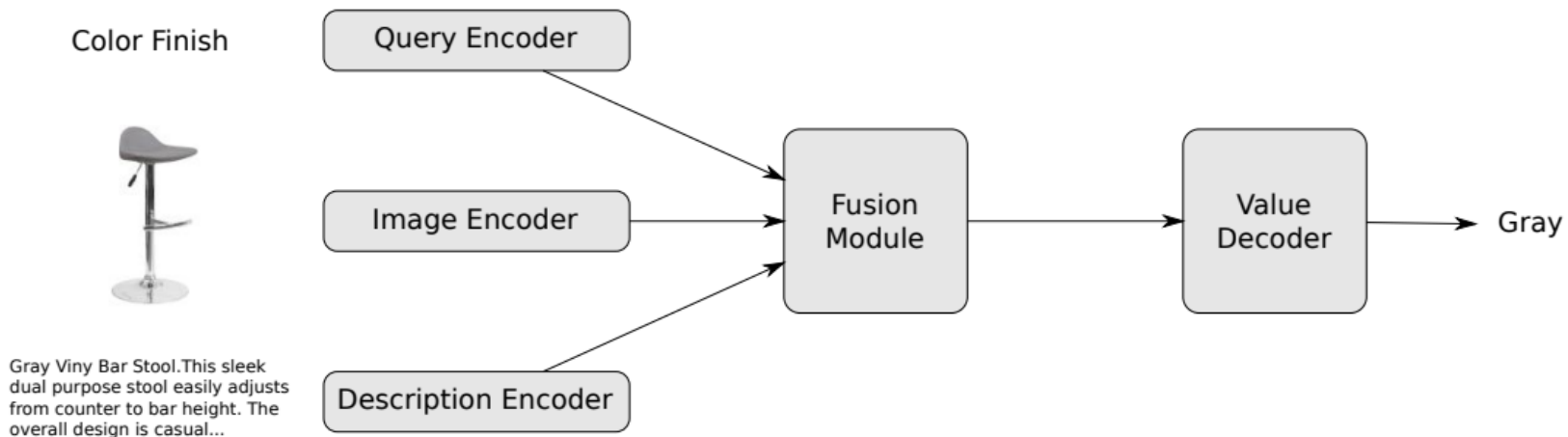


Messi celebrates with teammate **Suarez** and **Griezmann** after hat-trick against Atletico Madrid.



- Multimodality
- Multimodal KG Construction
- Inference
- Challenges

Multi-modal KG Completion



Logan IV R L, Humeau S, Singh S. Multimodal attribute extraction. NIPS, 2017.

Multi-modal KG Completion

Input

Item Details



Output

商品名称: 芝华仕沙发	商品编号: 100000993577	商品毛重: 116.0kg	商品产地: 重庆、苏州、惠州、...
货号: 8908A	功能类别: 功能沙发	材质类别: 布艺沙发	适用面积: 中户型 (60m ² -90m ²)
适用人数: 三人	填充物: 海绵	沙发组合样式: 一字型	是否可定制: 不可定制
面料材质: 科技布	是否带储物: 不可储物	主色系: 白色/浅色	是否带贵妃榻: 不带
风格: 现代简约	是否可拆洗: 不可拆洗	框架材质: 板木结合	适用场景: 客厅, 庭院, 阳台

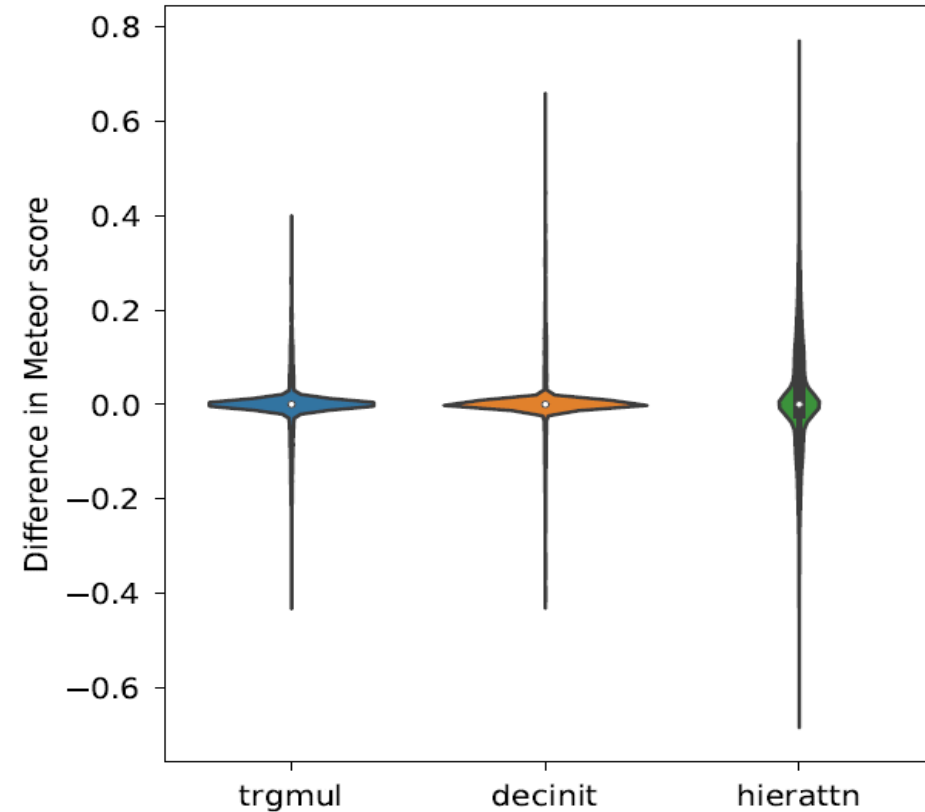
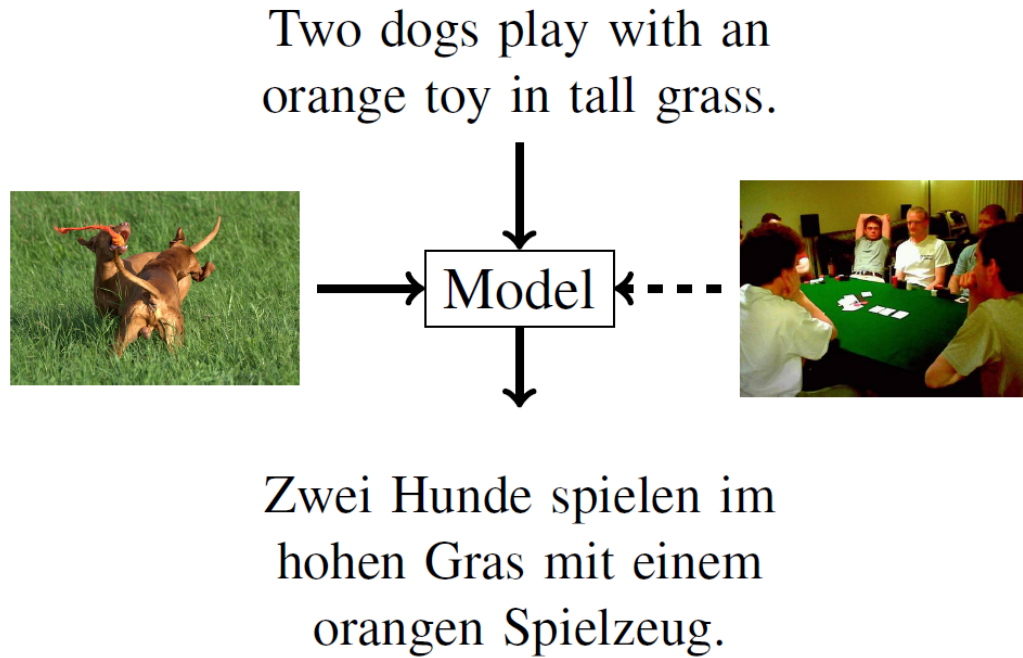
[三人、白色/浅色、不带]



Our MKG base model

Is the multi-modal KG really helpful?

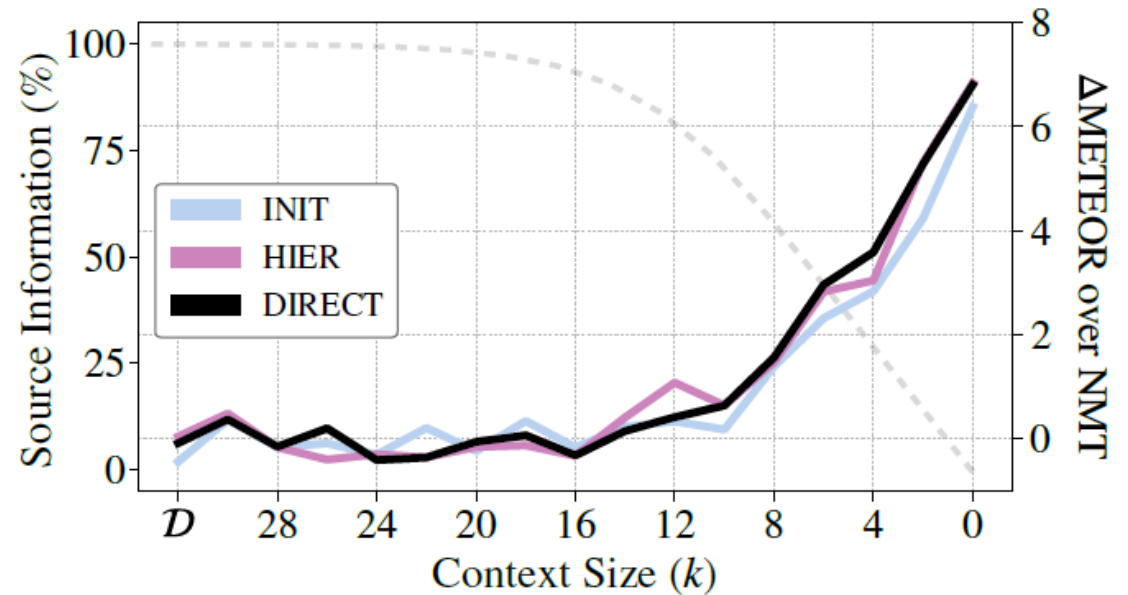
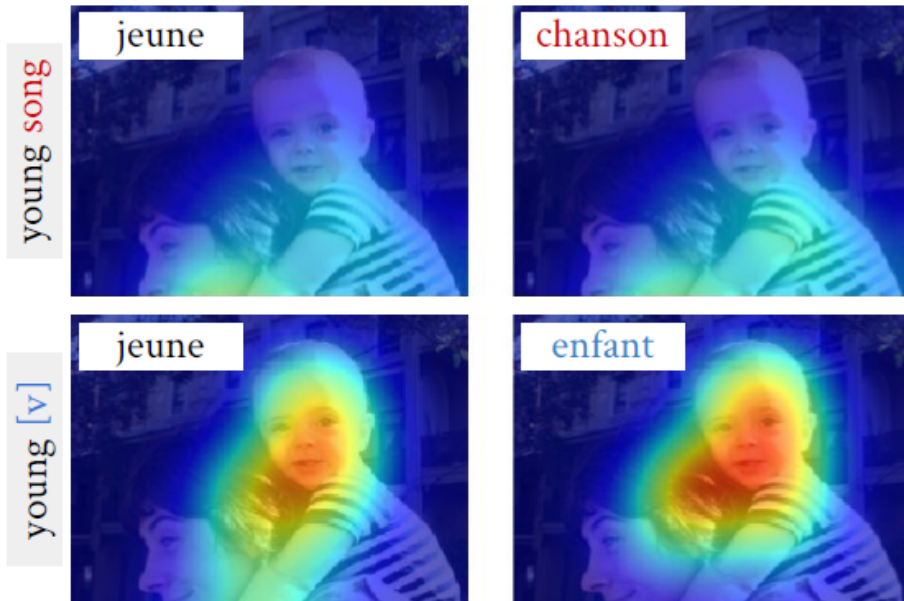
Adversarial Evaluation of Multimodal Machine Translation. EMNLP 2018. (CCF B)



Only needed for incorrect, ambiguous, and gender-neutral words

Is the multi-modal KG really helpful?

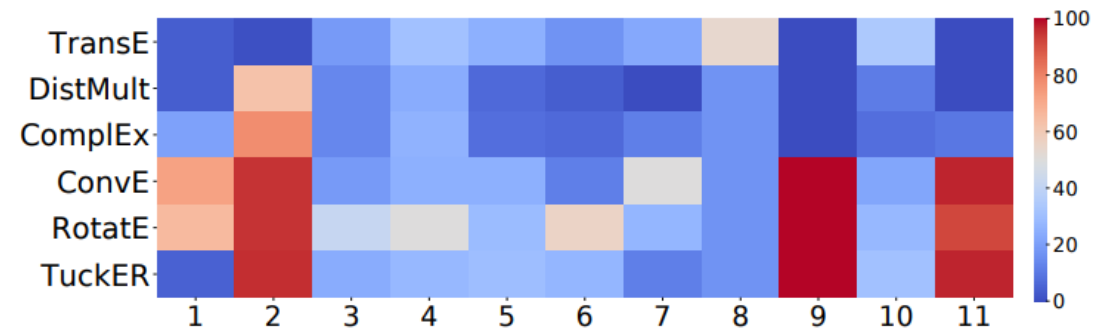
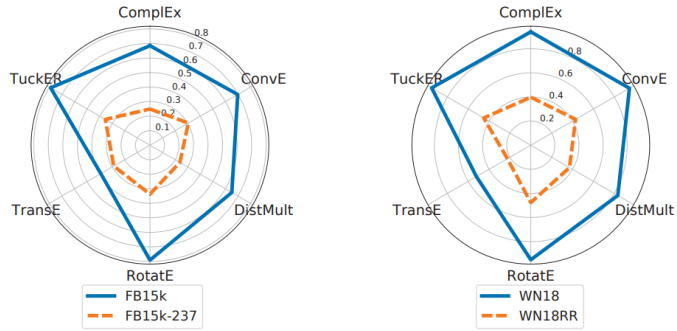
Probing the Need for Visual Context in Multimodal Machine Translation.
NAACL 2019. (CCF B)



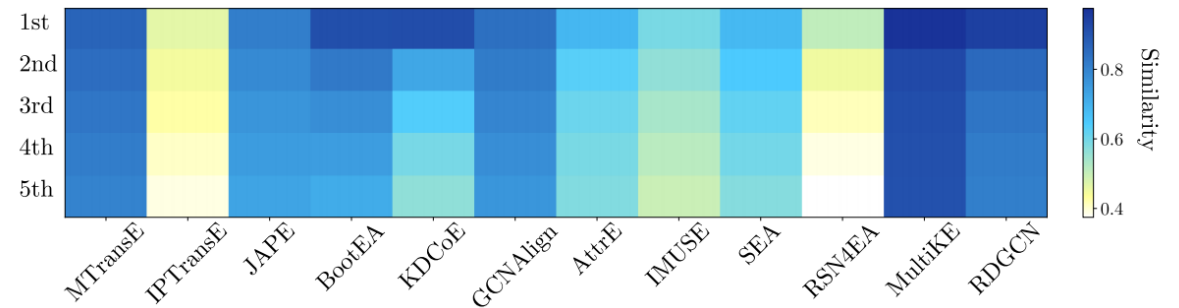
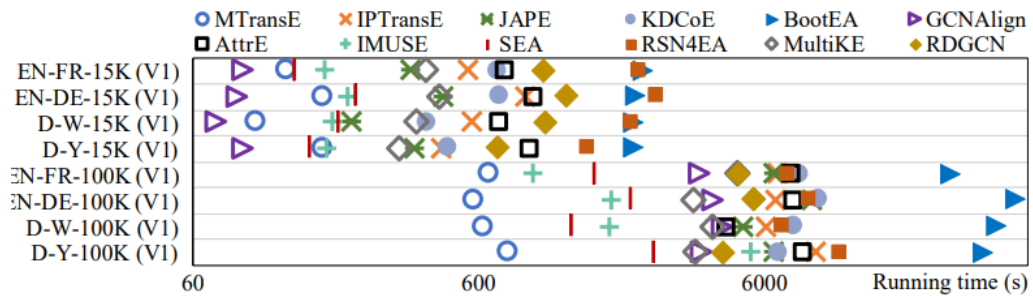
This dominance effect corroborates the seminal work of Colavita (1974) in Psychophysics where it has been demonstrated that visual stimuli dominate over the auditory stimuli when humans are asked to perform a simple audiovisual discrimination task.

Is the multi-modal KG really helpful?

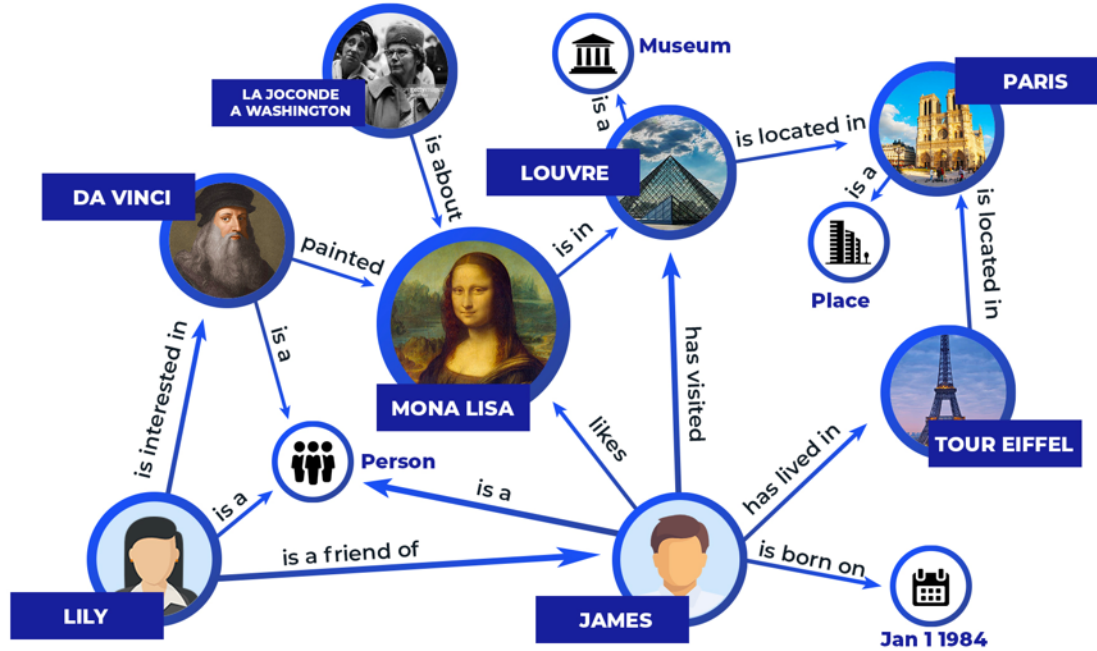
Realistic Re-evaluation of Knowledge Graph Completion Methods: An Experimental Study. SIGMOD 2020. (CCF A)



A Benchmarking Study of Embedding-based Entity Alignment for Knowledge Graphs. PVLDB 2020. (CCF A)



Is the multi-modal KG really helpful?



Cross-modal Entity Linking ? ? ?

Is the multi-modal information only needed in very specific cases for KG?

- Multimodality
- Multimodal KG Construction
- Inference
- Challenges

Challenges

Parsing text to structured semantic graph

Parsing images/videos to structures

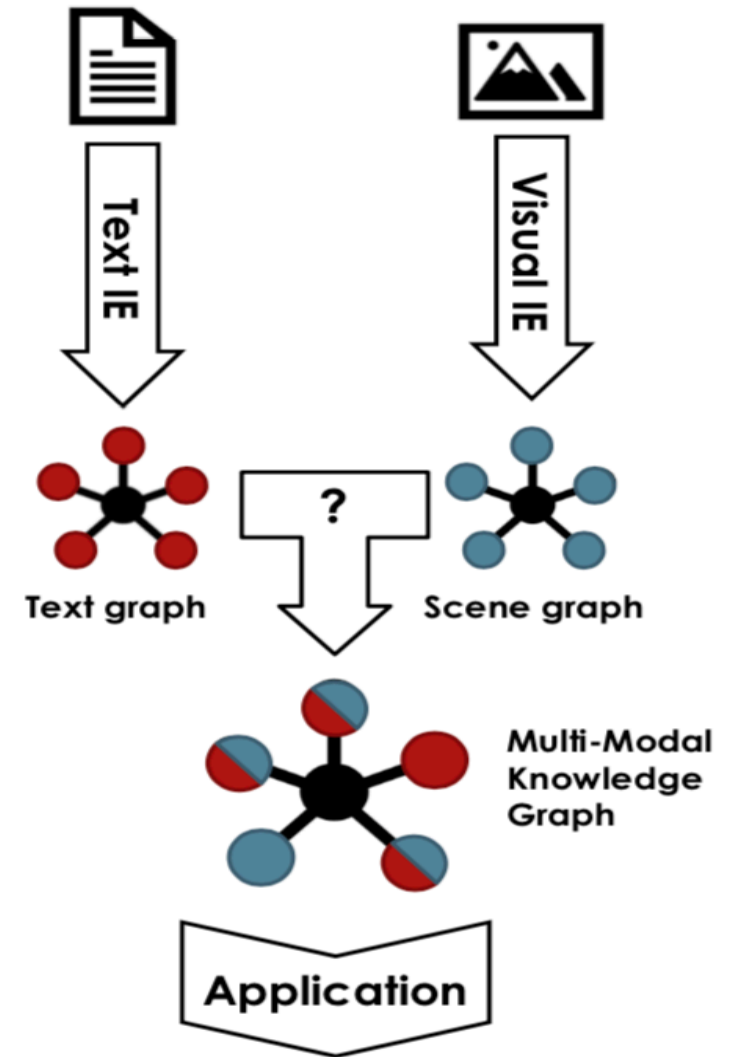
Grounding event/entities across modalities

Annotation Cost or Limited training data (domain specific)

Computational complexity

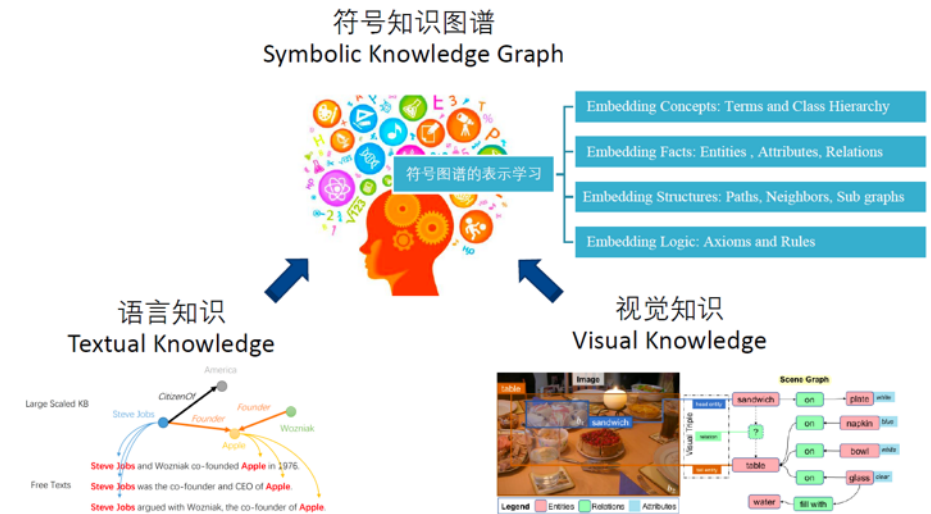
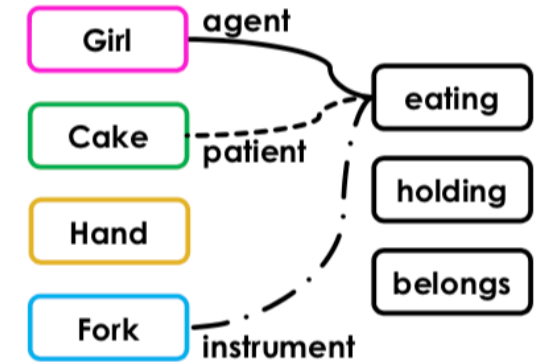
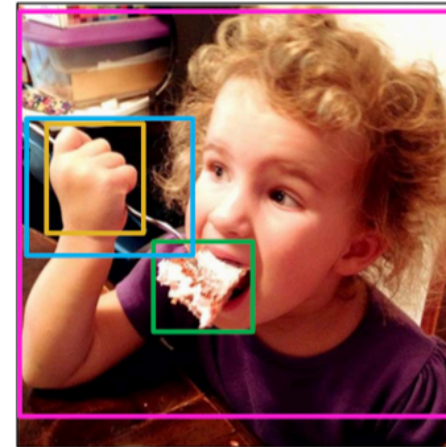
Limited fixed vocabulary

Abstract concept not groundable

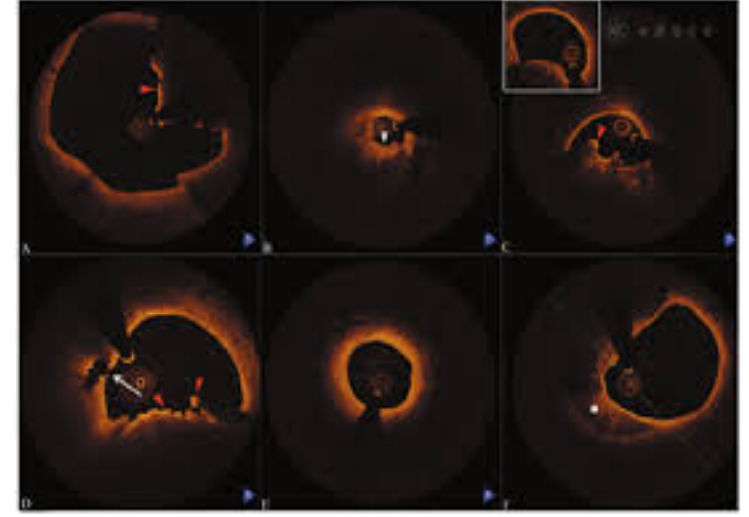
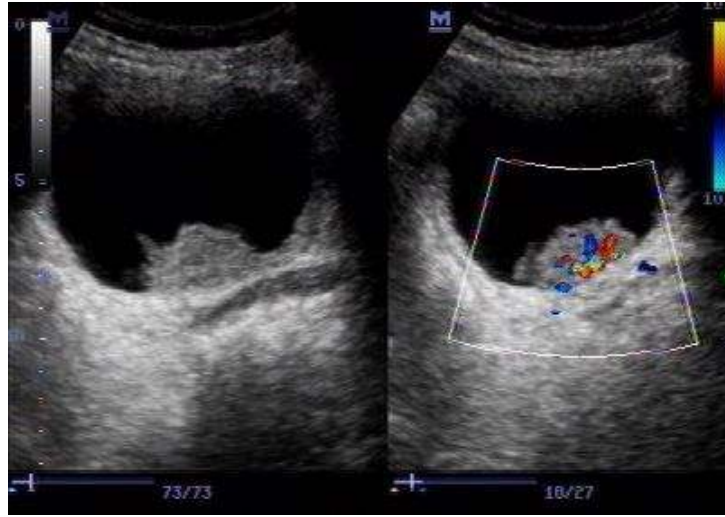
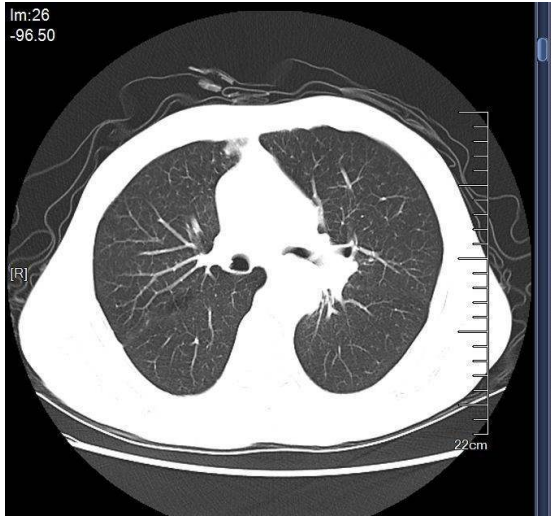


Real Challenges

- **Multimodal Data:**
 - **KG**
 - **Text**
 - **Image or video**
- **Multimodal Knowledge Representation:**
 - **Multimodal**
 - **Spatial-Temporal**
 - **Event**
 - **Rules**
- **Multimodal Representation Learning:**
 - **Pre-trained model for multimodal KG**
 - **Cross-modal alignment**
 - **Computing and storage capacity**

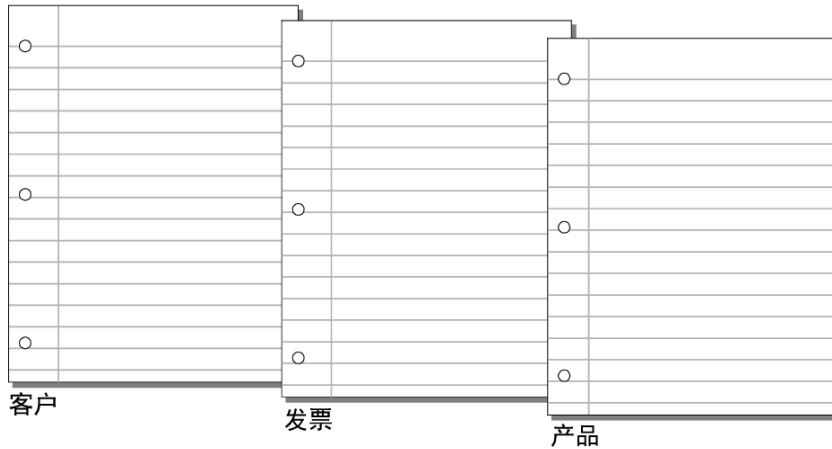


Multimodal Image?

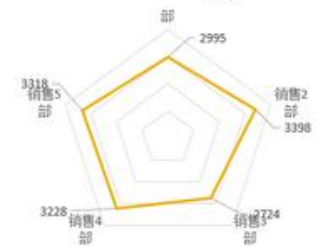


Other Multimodal Data?

Multimodal Tabular Data?



销售部	第一季度	第二季度	第三季度	第四季度	合计	辅助列	
						第一季度	产品
销售1部	699	630	892	774	2995	699	销售1部
销售2部	800	988	814	796	3398	800	销售2部
销售3部	870	590	646	618	2724	870	销售3部
销售4部	1098	488	992	650	3228	1098	销售4部
销售5部	1234	458	936	690	3318	1234	销售5部
合计	4701	3154	4280	3528	15663	10962	三季度合计



Feature	Elements
Visual	Composition, Texture, Size, Color, Saturation, Focus
Motion	Zooming/Tracking, Camera Position, PerspectiveSpeed, Pan/Tilt, Editing
Audio	Volume, Speed, Pitch, Music, Tone, Frequency
Text	Size, Placement, Color, Diction, Tone, Font

Richpedia Demo

Thank you!