# Multi-Modal Knowledge Graph Construction and Application: A Survey

Xiangru Zhu, Zhixu Li *Member, IEEE,* Xiaodan Wang, Xueyao Jiang, Penglei Sun, Xuwu Wang,
Yanghua Xiao *Member, IEEE,* Nicholas Jing Yuan *Member, IEEE,*

**Abstract**—Recent years have witnessed the resurgence of knowledge engineering which is featured by the fast growth of knowledge graphs. However, most of existing knowledge graphs are represented with pure symbols, which hurts the machine's capability to understand the real world. The multi-modalization of knowledge graphs is an inevitable key step towards the realization of human-level machine intelligence. The results of this endeavor are Multi-modal Knowledge Graphs (MMKGs). In this survey on MMKGs constructed by texts and images, we first give definitions of MMKGs, followed with the preliminaries on multi-modal tasks and techniques. We then systematically review the challenges, progresses and opportunities on the construction and application of MMKGs respectively, with detailed analyses of the strengths and weaknesses of different solutions. We finalize this survey with open research problems relevant to MMKGs.

**Index Terms**—Multimodal Knowledge Graph, Survey, Symbol Grounding

◆

## 1 INTRODUCTION

RECENT years have witnessed the resurgence of knowledge engineering featured by the fast growth of knowledge graphs. A knowledge graph (KG) is essentially a large-scale semantic network that contains entities, concepts as nodes and various semantic relationships among them as edges. The great value of knowledge graphs has been found in a wide range of real-world applications, including text understanding, recommendation systems and natural language question answering. More and more knowledge graphs have been created, covering common sense knowledge (e.g., Cyc [1], ConceptNet [2]), lexical knowledge (e.g., WordNet [3], BabelNet [4]), encyclopedia knowledge (e.g., Freebase [5], DBpedia [6], YAGO [7], WikiData [8], CN-DBpedia [9]), taxonomic knowledge (e.g., Probase [10]) and geographic knowledge (e.g., GeoNames [11]).

However, most of the existing knowledge graphs are represented with pure symbols denoted in the form of text, which weakens the capability of machines to describe and understand the real world. A human being cannot understand what a `dog` is without the experience of living with a dog, which enlightens researchers to establish the connection between the symbol `Dog` and the experience of dogs, that is, grounding a symbol to its physical world meaning [12], [13], [14]. Similarly, grounding symbolic forms to non-symbolic experiences benefits receiving real communicative intents [15]. For example, the customers cannot understand the meaning of `Hand-in-waistcoat` as a particular pose (hand inside coat flap) without the *experience* of `Hand-in-waistcoat` so that the customer would respond incorrectly to the request of photographers. Thus, it is nec-

essary to ground symbols to corresponding images, sound and video data and map symbols to their corresponding referents with meanings in the physical world, enabling machines to generate similar "*experiences*" like a real human [12] when they are confronted with a specific entity `Hand-in-waistcoat` or an abstract concept `Dog`. On the other hand, there is an increasing demand for the multi-modality of knowledge to break through the bottleneck of real-world applications [16], [17], [18]. For instance, in relation extraction tasks, an additional image usually greatly improves the performance in the extraction of the attributes and relationships that are visually obvious but difficult to be recognized in symbols and text, such as *partOf* (e.g., *The keyboard and the screen are parts of a laptop.*) and *colorOf* (e.g., *A banana is usually yellow or yellowish-green but not blue* ). In text generation tasks, if the machine has been empowered with the ability to recognize a specific entity in an image by the reference to a Multi-Modal KG (MMKG), the machine is possible to generate a more informative entity-level sentence (e.g., *Donald Trump is making a speech*) instead of a vague concept-level description (e.g., *A tall man with blond hair is making a speech*).

Due to the rapid growth of applications' demand for multi-modal knowledge guidance, the multi-modalization of KGs and their applications has been booming in recent years [19], [20], [21]. Nevertheless, a systematic review of the recent research progresses, challenges and opportunities in this emerging area are still lacking. In this paper, we hope to fill the gap and systematically survey the recent research progresses relevant to MMKG as follows: 1) **Construction.** The construction of MMKGs could be conducted in two opposite directions. One is from images to symbols, i.e., labeling images with symbols in KG; the other is from symbols to images, i.e., grounding symbols in KG to images. In the Construction section, we will systematically cover the challenges, progresses as well as opportunities to correlate various symbol knowledge (e.g., entities, concepts, relations

---

- X. Zhu, Z. Li, X. Wang, X. Jiang, P. Sun, X. Wang and Y. Xiao are with the School of Computer Science, Fudan University.
  E-mail: {xrzhu19, zhixuli, xiaodanwang20, xueyaojiang19, plsun20, xwwang18, shawyh}@fudan.edu.cn. Z. Li and Y. Xiao are the corresponding authors.
- N.J. Yuan is with Huawei Cloud & AI, Hangzhou, Zhejiang, China.
  E-mail: nicholas.yuan@huawei.com

(a) MMKG with multi-modal data as attribute values
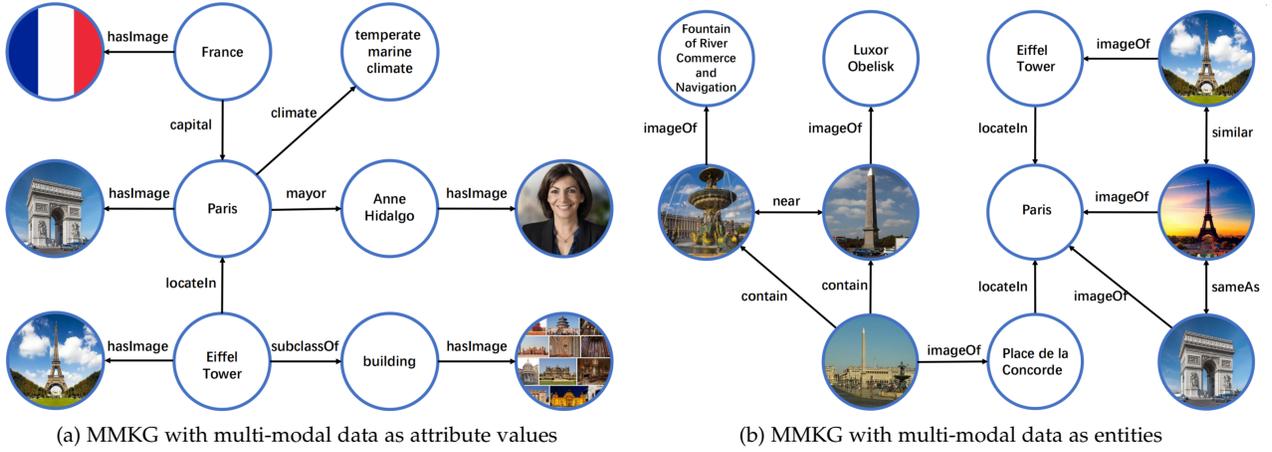


(b) MMKG with multi-modal data as entities

Fig. 1: Example MMKGs of two different types: A-MMKG and N-MMKG

and events) to their corresponding images in the two opposite directions. 2) **Application.** The application of MMKGs could be roughly divided into two categories: In-MMKG applications aiming at addressing the quality or integration issues of MMKGs themselves, and Out-of-MMKG applications which are general multi-modal tasks that MMKGs can help. The Application section will present how MMKGs are applied in several well-studied multi-modal tasks.

To summarize, we are the first to thoroughly survey the existing work on MMKGs consisting of texts and images. To enhance the value of this survey, we pay attention to the following features: 1) **Comprehensive Survey.** We systematically and comprehensively review the existing work on MMKG construction and application. 2) **Insightful Analysis.** We analyze the strengths and weaknesses of different solutions in MMKG construction and discuss how MMKGs can help in various downstream applications. 3) **Revealed Opportunities.** We not only point out some potential opportunities with the studied tasks relevant to MMKG construction, but also list some promising future directions with MMKG.

The rest of the survey is organized as follows: Sec. 2 gives definitions and preliminaries on MMKGs. Sec. 3 conducts a comprehensive review of the challenges, progresses and opportunities of the construction of MMKGs, while Sec. 4 presents how MMKGs are applied in several well-studied multi-modal applications. Sec. 5 reviews some open problems of MMKG and highlights promising future directions. Sec. 6 finally concludes the paper.

## 2 DEFINITIONS AND PRELIMINARIES

This section first defines two representation ways for KGs and then reviews some preliminaries on multi-modal techniques, followed by a discussion on the connections between MMKGs and the existing multi-modal techniques.

### 2.1 Definition & Representation of MMKGs

A traditional Knowledge Graph (KG) is defined as a directed graph $\mathcal{G} = \{\mathcal{E}, \mathcal{R}, \mathcal{A}, \mathcal{V}, \mathcal{T}_{\mathcal{R}}, \mathcal{T}_{\mathcal{A}}\}$, where $\mathcal{E}$, $\mathcal{R}$, $\mathcal{A}$, $\mathcal{V}$ are sets of entities, relations, attributes and literal attribute values, and $\mathcal{T}_{\mathcal{R}} = \mathcal{E} \times \mathcal{R} \times \mathcal{E}$ and $\mathcal{T}_{\mathcal{A}} = \mathcal{E} \times \mathcal{A} \times \mathcal{V}$ are sets of relation triples and attribute triples respectively. A triple $(s, p, o) \in \mathcal{T}_{\mathcal{R}}$ denotes that *entity* $s \in \mathcal{E}$ has a *relation* $p \in \mathcal{R}$ with *entity* $o \in \mathcal{E}$. A triple $(s, p, o) \in \mathcal{T}_{\mathcal{A}}$ denotes that *entity* $s \in \mathcal{E}$ has an *attribute* $p \in \mathcal{A}$ with the *attribute value* $o \in \mathcal{V}$.

A Multi-modal Knowledge Graph (MMKG) can be seen as a multi-modalized KG, which has part of its knowledge in $\{\mathcal{E}, \mathcal{R}, \mathcal{A}, \mathcal{V}, \mathcal{T}_{\mathcal{R}}, \mathcal{T}_{\mathcal{A}}\}$ multi-modalized. We say a particular knowledge symbol is multi-modalized if it is associated with its corresponding data items in modalities other than text, such as image, sound or video, that could embody the knowledge. For instance, a relation triple $(s, p, o)$ can be multi-modalized with an image describing the *relation* $p$ between $s$ and $o$.

Existing work on MMKGs mainly adopts two different ways for representing MMKGs. One way takes multi-modal data (images in this survey) as particular attribute values of entities or concepts, as the example shown in Fig. 1(a). We name an MMKG represented in this way as **A-MMKG** for short, denoted as $\mathcal{G} = \{\mathcal{E}, \mathcal{R}, \mathcal{A}, \mathcal{V}, \mathcal{T}_{\mathcal{R}}, \mathcal{T}_{\mathcal{A}}\}$, where $\mathcal{T}_{\mathcal{A}} = \mathcal{E} \times \mathcal{A} \times (\mathcal{V}_{\mathcal{KG}} \cup \mathcal{V}_{\mathcal{MM}})$ is the set of attribute triples, $\mathcal{V}_{\mathcal{KG}}$ is the set of the KG's attribute values and $\mathcal{V}_{\mathcal{MM}}$ is the set of multi-modal data. In A-MMKGs, since multi-modal data are treated as attribute values, in a triple $(s,p,o)$, $s$ denotes an entity, $o$ denotes one of its corresponding multi-modal data, and the relation $p$ is "hasImage" when $o$ is an image. Some example triples are listed in Table 1a.

The other way takes multi-modal data as entities in KGs, as the example shown in Figure 1(b). We name an MMKG represented in this way as **N-MMKG** for short, denoted as $\mathcal{G} = \{\mathcal{E}, \mathcal{R}, \mathcal{A}, \mathcal{V}, \mathcal{T}_{\mathcal{R}}, \mathcal{T}_{\mathcal{A}}\}$, where $\mathcal{T}_{\mathcal{R}} = (\mathcal{E}_{\mathcal{KG}} \cup \mathcal{E}_{\mathcal{MM}}) \times \mathcal{R} \times (\mathcal{E}_{\mathcal{KG}} \cup \mathcal{E}_{\mathcal{MM}})$ is the set of relation triples, $\mathcal{E}_{\mathcal{KG}}$ is the set of KG entities and $\mathcal{E}_{\mathcal{MM}}$ is the set of multi-modal data. Since multi-modal data are treated as new entities, more inter-modal and intra-modal relations are discovered and added into the MMKG. For example, in Table 1b, the entity `Eiffel Tower` is associated with an image *Eiffel_Tower.jpg* by the relation `imageOf`. Two images can also be associated in one of the following relations: 1) `contain`: One image entity visually contains another image entity by the relative position of images. 2) `nearBy`: One image entity is visually nearby another image entity in an image. 3) `sameAs`: Two different image entities refer to the

| subject | predicate | object |
|---|---|---|
| France | hasImage | The_flag_of_France.jpg |
| Anne Hidalgo | hasImage | Anne_Hidalgo.jpg |
| Paris | mayor | Anne Hidalgo |
| Paris | hasImage | A_landmark_of_Paris.jpg |
| Eiffel Tower | locateIn | Paris |
| Eiffel Tower | hasImage | Eiffel_Tower.jpg |
| Eiffel Tower | subclassOf | building |
| building | hasImage | a_kind_of_architectural_style.jpg |

(a) Example RDF triples in A-MMKG

| subject | predicate | object |
|---|---|---|
| Eiffel_Tower_in_Paris.jpg | imageOf | Paris |
| Eiffel_Tower_in_Paris.jpg | size | 700*1600 |
| Eiffel_Tower_in_Paris.jpg | sameAs | Arc_de_Triomphe_in_Paris.jpg |
| Eiffel_Tower_in_Paris.jpg | similar | Eiffel_Tower.jpg |
| Eiffel Tower | locateIn | Paris |
| Eiffel_Tower.jpg | imageOf | Eiffel Tower |
| Eiffel_Tower.jpg.HOG | describes | Eiffel_Tower.jpg |
| Eiffel_Tower.jpg.HOG | value | [0.0775 , 0.0120 , 0.0021 , ...] |

(b) Example RDF triples in N-MMKG

TABLE 1: Example RDF triples in different types of MMKGs, where items end up with ".jpg" are images.

same entity. 4) `similar`: Two image entities are visually similar to each other.

In addition, in N-MMKGs an image is usually abstracted into several image descriptors, which are usually summarized into feature vectors of the image entity at the pixel level, such as Gray Histogram Descriptor (GHD), Histogram of Oriented Gradients Descriptor (HOG), Color Layout Descriptor (CLD) and so on. For example, in Table 1b, *Eiffel_Tower_in_Paris.jpg.HOG* is one of the descriptors of the image *Eiffel_Tower_in_Paris.jpg*, and is in the form of a vector. These image descriptors are well interpreted. Thus the relations between images can be obtained by simple calculations (e.g., image similarity obtained via the inner product of vectors of image descriptors).

We list mainstream MMKGs constructed with image-based visual knowledge extraction systems in Table 2(a). NEIL [22] annotates each image with a single label by pre-trained classifiers and extracts visual relations by heuristic rules about the locations of extracted objects. GAIA [21] extracts fine-grained concepts in the news by object recognition together with fine-grained classification. Based on the framework of GAIA, RESIN [23] extracts visual news events and identifies related visual entities and concepts as arguments on small-scale resources (news documents). Later, MMEKG [24] optimizes some modules and adapt to billion-scale universal events extraction.

MMKGs listed in Table 2(b) are constructed with symbol grounding. IMGpedia [25] linkes images from Wikimedia Commons[1] to DBpedia via the structured Wikipedia data in RDF format already extracted in DBpedia Commons [26], which additionally adds the similarity relations between images. ImageGraph [27] searches images from search engines with entities in KGs as queries. MMKG [28] extends this method to several KGs by aligning entities across them. These construction methods based on symbolic entity alignment (such as by linked datasets or URI of entities) focus on the representativeness of images, but the diversity of images is also an important issue due to different contexts and views. Richpedia [29] trains an additional diversity retrieval model to select diverse images. The categories of entities in Richpedia are limited to cities, sights, and persons. In addition, VisualSem [30] considers that many entities are non-visualizable entities that should not be searched for images. Therefore it starts with the most typical visual entities and mine other related visual entities iteratively. However, the small scale of VisualSem is far from satisfying the knowledge demands of downstream applications.

1. a multi-media dataset linking to Wikipedia articles, https://wikimediafoundation.org/our-work/commons/

## 2.2 Preliminaries on Multi-Modal Techniques

Modality refers to the particular way in which something exists, is experienced or is done [31]. In computer science and artificial intelligence, a problem is characterized as multi-modal if it involves data of multiple modalities. Typical multi-modal tasks with images and texts include image caption [32], visual question answering [33], and cross-modal retrieval [34], etc. We will introduce how MMKGs are applied in these applications in Sec. 4.2. But before MMKGs, people mainly focus on multi-modal learning, and more recently the Vision and Language Pre-trained Models (VL-PTMs), which will be briefly introduced below.

**Multi-Modal Learning**. Multi-modal learning focuses on modeling the correspondences among multiple modalities, which includes: 1) *Multi-modal Representation* aims to use the complementary of multi-modality to learn feature representation. The existing efforts either project the multiple modalities into a unified space [35], or represent every single modal in its own vector space which satisfies certain constraints like linear correlation [36]. 2) *Multi-modal Translation* learns to translate from a source instance in one modality to a target instance in another, including example-based [34] and generative translation models [37]. 3) *Multi-modal Alignment* aims to find the correspondences between different modalities. It can either be directly applied in some multi-modal tasks such as visual grounding or as a pre-training task in VL-PTMs [38]. 4) *Multi-modal Fusion* aims to join information from different modalities to perform a prediction [31], where various attention mechanisms [39], [40] are applied to model the interaction between different features in the cross-modal module. 5) *Multi-modal Co-Learning* aims to alleviate the low-resource problems in a certain modality by leveraging the resources of other modalities through the alignment between them [31].

**VL-PTMs**. Recently, many large companies and research institutions including OpenMind [41], Microsoft [42], [43] and Huawei [44] etc. pay great efforts on training large VL-PTMs based on large-scale unsupervised multi-modal data. A typical VL-PTM example is CLIP [41] trained on 400 million text-image pairs, which significantly improves the performance of image classification and cross-modal retrieval. Based on massive multi-modal data and large-scale models, VL-PTMs could learn extensive implicit cross-modal knowledge with some designed self-supervised pre-training tasks, such as masked language model, sentence image alignment, masked region label classification, masked region features regression, masked object prediction, etc. Furthermore, to improve fine-grained cross-modal understanding, some work also add cross-modal object align-

| System | MMKG Type | Multi-modalized Knowledge | Source Images | Candidate KGs | Quality Control | Scale |
|---|---|---|---|---|---|---|
| NEIL [22] | N-MMKG | entity, concept, relation | images from search engine | WordNet | semi-supervised classification with labeled seed images | 1,152 objects,1,034 scenes 87 attributes,1,703 triples (2.5 months) |
| GAIA [21] | N-MMKG | entity, concept | multimedia news documents | Freebase, GeoNames | object detection, fine-grained classification, heuristic rules | < 457K entities, < 67K triples, < 38K events (including textual and visual ones) |
| RESIN [23] | N-MMKG | entity, concept, event | multimedia news documents | WikiData | event classification, object recognition, situation recognition, weakly-supervised event grounding, event relation extraction | < 24 entities, < 46 relations, <67 events (including textual and visual ones) |
| MMEKG [24] | N-MMKG | event | Wikipedia, BookCorpus, CC3M&CC12M, C4(news) | WordNet | event classification, object recognition, event relation extraction | < 990K events, < 644 event relations < 863M instance events, < 934M instance events' relations (including textual and visual ones) |

(a) Image-based visual knowledge extraction systems that could be used to construct MMKGs by labeling images

| MMKG | MMKG Type | Multi-modalized Knowledge | Source KGs | Candidate images | Quality Control | Scale | Images per entity |
|---|---|---|---|---|---|---|---|
| IMGpedia [25] | N-MMKG | entity, concept, relation | DBpedia | Wikimedia Commons | constructed via DBpedia Commons | 12.7M links to KG (with 2.6M DBpedia entities /concepts), 3000M triples (including 443M triples of 1 visual relation) | >5.6 |
| ImageGraph [27] | A-MMKG | entity, concept | FB15K | search engine | disambiguation by Wikipedia URI | 15K entities/concepts | 55.8 |
| MMKG [28] | A-MMKG | entity, concept | FB15K, DBpedia15K, YAGO15K | search engine | 1.entity alignment cross different KGs 2.disambiguation by Wikipedia URI | 15K entities/concepts | 55.8 |
| Richpedia [29] | N-MMKG | entity, concept, relation | Wikidata | search engine, Wikipedia | 1.disambiguation by Wikipedia URI 2.a diversity retrieval model to filter images | 2.8M entities/concepts, 172M triples (including 114.5M triples of 3 visual relations) | 99.2 |
| VisualSem [30] | N-MMKG | entity, concept | BabelNet | Wikipedia, ImageNet | 1.synsets in ImageNet as initial entities pool 2.mining neighbours 3.a image-text matching model to filter noise | 89.9K entities/concepts, 13 relations | 10.4 |

(b) MMKGs constructed by symbol grounding

TABLE 2: Mainstream MMKGs (or extraction systems for constructing MMKGs) and their relevant information

ment [43], [44], [45], relation alignment [46], [47] tasks to optimize the pre-training process.

## 2.3 Discussions

Although there there is already much research on multi-modal learning and VL-PTMs, there is still an emerging trend to introduce MMKGs to help enhance multi-modal tasks. In general, MMKGs could benefit multi-modal tasks in the following aspects.

1) MMKGs provide sufficient background knowledge to enrich the representation of entities and concepts, especially for the long-tail ones. For instance, [16] uses auxiliary commonsense knowledge to enhance the representation of image and text to improve image-text matching.
2) MMKGs enable the understanding of unseen objects in images. Unseen objects pose a great challenge to statistic-based models. Symbolic knowledge alleviates the difficulty by providing symbolic information about unseen objects or establishing semantic relations between seen objects and unseen objects. For example, [48] uses external symbolic knowledge to guide the generation of captions for unseen novel visual objects.

3) MMKGs enable multi-modal interpretable reasoning. For example, the OK-VQA dataset [49], which contains only questions that require external knowledge to answer, is built to test the reasoning capability of VQA models.
4) MMKGs usually provide multi-modal data as additional features to bridge the information gaps in some NLP tasks. In the case of entity recognition, the image could provide sufficient information to identify whether "Rocky" is the name of a dog or a person [50].
5) MMKGs provide explicit and fine-grained cross-modal correlation knowledge, which is complementary to the implicit knowledge learned by VL-PTMs. Besides, MMKGs have advantages on providing long-tail knowledge, background knowledge, and fine-grained knowledge compared with VL-PTMs [51].

To sum up, previous efforts to use multi-modal information are still limited without the support of large-scale MMKG. Multi-modal tasks can be further improved when MMKGs are available.

## 3 CONSTRUCTION

The essence of MMKG construction is associating symbolic knowledge in a traditional KG, including entities, concepts,
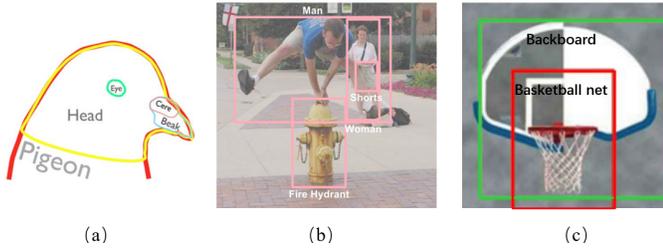
Fig. 2: Examples of labeling images: (a) labeling components after image segmentation in Visipedia [52]; (b) labeling objects with bounding boxes in Visual Genome [53]; (c) labeling two objects where one is a part of the other in NEIL [22], e.g., *PartOf*(`Basketball net`, `Backboard`).

relations, etc., with their corresponding images. Two opposite ways to complete the task are (1) labeling images with symbols in KG and (2) grounding symbols in KG to images. We elaborate on the two categories of solutions in Sec. 3.1 and Sec. 3.2 respectively. We finally discuss the differences between the two solutions in Sec. 3.3.

## 3.1 From Images to Symbols: Labeling Images

The CV community has developed many image labeling solutions, which could be leveraged in labeling images with structural symbols (e.g., concepts or entities) in KG. For example, NEIL [22] links images to WordNet [3], and ImageSnippets [54], [55] links images to DBPedia [6]. Most image labeling solutions learn the mapping from image content to a wide variety of label sets, including objects, scenes, entities, attributes, relations, events and other symbols. The learning procedure is supervised by human-annotated datasets, which require the crowd workers to draw bounding boxes and annotate images or regions of images with given labels, as illustrated in Figure 2.

Some well-known image-based visual knowledge extraction systems are as listed in Table 2(a), which could be utilized for constructing MMKGs through image labeling. According to the category of symbols to be linked, the process of linking images to symbols could be divided into several fractionized tasks: *visual entity/concept extraction* (Sec. 3.1.1), *visual relation extraction* (Sec. 3.1.2) and *visual event extraction* (Sec. 3.1.3).

### 3.1.1 Visual Entity/Concept Extraction

Visual entity (or concept) extraction aims to detect and locate target visual objects in images and then label these objects with entity (or concept) symbols in KG.

**CHALLENGES**. The main challenge with this task lies in how to learn an effective fine-grained extraction model without a large-scale, fine-grained, well-annotated concept and entity image dataset. Although there are rich well-annotated image data in CV, these datasets are almost coarse-grained concept images, which could not meet the requirements of MMKG construction for image annotation data of fine-grained concepts and entities.

**PROGRESSES**. The existing efforts with visual entity/concept extraction could be roughly divided into two categories: 1) object recognition methods, which label a



Fig. 3: The heatmap for detected visual entities (`Soldier` and `Boats`) in two example images by visual grounding in GAIA [21], where the stronger the correlation between a pixel and a word, the warmer the color of the pixel.
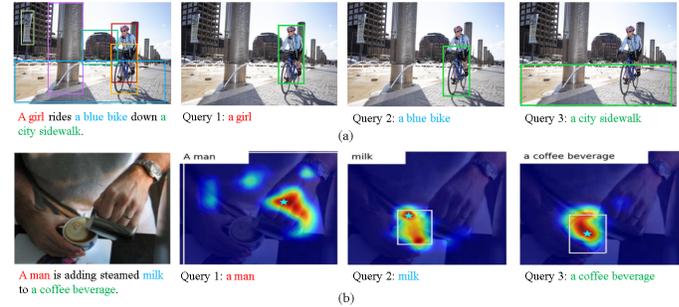


Fig. 4: Two kinds of weakly supervised visual entity extraction: (a) the attention-based method [56] and (b) the saliency-based method [57]. The first method selects the most relevant bounding boxes to given phrases. The second method selects the most sensitive pixels to given phrases.

visual entity/concept by classifying the region of a detected object; and 2) visual grounding methods, which label a visual entity/concept by mapping a word or phrase in a caption to the most relevant region.

1) *Object Recognition Methods.* In early works, images provided by users and researchers are usually simple and there is only one object in one image, which can be processed by classification models. But images in our real life could be too complex to be represented with only one label. Thus we need to tag different visual units with different labels.

In order to distinguish several visual entities in images, pre-trained detectors and classifiers are needed to label visual entities (as well as attributes and scenes) with their locations in the images. These detectors are trained with supervised data from public images-text datasets [21] (such as MSCOCO [58], Flickr30k [59], Flick30k Entities [60] and Open Images [61]). In the detection process, detectors (e.g., face detectors based on MTCNN or vehicle detectors based on Faster-RCNN) capture a set of region proposals for possible objects. In the recognition process, the pre-trained classifiers pick out region proposals that do contain objects and recognize candidate visual objects with entity-level (e.g., `BMW 320`) or concept-level (e.g., `Car`) labels. Since many recognized objects are duplicated instances of the same entities at different viewpoints, positions, poses and appearances, a common way to process is to cluster all the regions with recognized objects, and only the central one of each cluster will eventually be the output as a new visual entity [21]. However, the disadvantage of these supervised solutions is that only a limited number of visual entities

under pre-defined labels could be recognized. The precision of the visual object extraction model used in GAIA is only 43% on the benchmark MSCOCO [21]. It requires much pre-processing work for fine-grained recognition [21], such as pre-defined rules, pre-trained fine-grained detectors, etc.

2) *Visual Grounding Methods.* In visual entity extraction, training detectors need a large amount of labeled data with bounding boxes and pre-defined schemas with a fixed set of concepts [62], which is challenging for large-scale visual knowledge acquisition. Fortunately, many image-captions pairs from the web weakly supervise the extraction of visual knowledge without relying on the labeled bounding boxes. Therefore, the visual entity extraction problem is reduced to an open-domain visual grounding problem, which aims to locate the corresponding image region of each phrase in a caption to obtain visual objects with their labels.

In the extraction process, we often select active pixels for the given word as the region of visual objects based on the spatial heatmap, as shown in Figure 3. In the cross-modal unified vector space, the heatmap of each phrase can be learned by attention-based methods and saliency-based methods, as shown in Figure 4. Saliency-based methods treat the marginal effects [63] of pixels to a given phrase by gradient computation as the heatmap value, and attention-based methods treat the cross-modal relevance as the heatmap value. However, since some salience methods are too sensitive to input changes to produce reliable results [64], [65], [66], thus attention-based methods [21], [56], [62], [67], [68] are more studied than saliency-based methods [57], [69] on locating visual objects. For example, the heatmap values in GAIA [21] are similarities between image regions and entity mentions in a caption, and those in [68] are similarities between image regions and possible event argument role types. At test time, the heatmap is thresholded to obtain a suitable bounding box of a visual object. If there is no overlap between the new bounding box and existing visual entities/concepts in KGs, the bounding box will be created as a new visual entity or concept.

The located visual objects via visual grounding include entities, concepts and attributes with acceptable accuracy. The accuracy of visual grounding methods used in GAIA [21] is 69.2% on Flickr30k. However, inconsistent semantic scales of images and texts may lead to incorrect matching. For example, `troops` may be mapped to *several individuals wearing military uniforms*, and `Ukraine (country)` may be mapped to *a Ukrainian flag*, both of which are relevant but not equivalent.

**OPPORTUNITIES**. 1) *VL-PTMs Based Extraction.* VL-PTMs bring new opportunities to nearly all cross-modal downstream tasks, including the detection of visual entities and concepts [70], [71]. The mapping of image patches and words can be directly visualized in the self-attention maps of the model without additional training. An example of the prediction with ViLT [72] is shown in Fig. 5. It is proved that VL-PTMs such as CLIP [41], trained on hundreds of millions of image-text data, can recognize many popular entities such as famous people and landmarks with high accuracy [73]. 2) *Taxonomy Extension.* Some visual objects with multiple reasonable labels indicate different semantic levels. For example, an image of *a boy* can be labeled



Fig. 5: Weakly supervised visual entity extraction via VL-PTMs. This figure shows the most relevant regions of an image to given words in a caption through self-attention mechanism of ViLT [72].

as `Person`, `Man` and `Boy`. To reduce the ambiguity, we should find an appropriate extension semantic level for the labels of images in the taxonomy. [53] fuses aforementioned multiple labels into the lowest common ancestor node of these synsets (i.e., `Person`), which may lead to many coarse-grained labels. [74] limits the scale of independent concepts' labels by setting a small value of maximum extension level to avoid too many related images. More nodes should be further searched recursively in the taxonomy consisting of hyponyms of the ancestor node to select the most semantically consistent label with the given visual object.

### 3.1.2 Visual Relation Extraction

Visual relation extraction aims to identify semantic relations among detected visual entities (or concepts) in images and then label them with the relations in KGs [22].

**CHALLENGES**. Although visual relation detection has been studied extensively in the CV community, most detected relations are superficial visual relationships between visual objects such as (`Person`, `standing on`, `Beach`). Differently, for the purpose of constructing MMKG, the visual relation extraction task aims to identify more general types of semantic relations that are defined in KGs, such as (`Jack`, `spouse`, `Rose`).

**PROGRESSES**. The existing efforts on visual relation extraction can be roughly put into two categories: rule-based relation extraction and statistic-based relation extraction. Some other work mainly focuses on long-tail relations and fine-grained relations, which will also be covered in the following.

1) *Rule-based Relation Extraction.* Traditional rule-based methods mainly focus on specific relations types, such as spatial relation [75], [76] and action relation [77], [78], [79], [80]. Experts usually predefine the criteria, and the discriminative features are scored and selected by heuristic methods.

In rule-based methods, the relations are determined based on label types of visual objects and the relative locations of regions. For example, if the bounding box of one object is always within that of another, there may be a `PartOf` relation between them. Table 3 lists several visual relations detected in NEIL, where the average detection accuracy of all 1703 relations is 79% [22]. During the extraction in NEIL, the detected relation between a pair of objects is, in turn, an additional constraint for new instance labeling. For example, "*Wheel is a part of Car*" indicates that it is more likely for a `Wheel` to appear in the bounding box of a `Car`. Rule-based methods provide highly accurate visual relations, but require much manual manipulation, which is less practical in large-scale MMKG construction.

| relation type | example | images | relation type | example | images |
|---|---|---|---|---|---|
| Concept-Concept | `Keyboard` is a part of `Laptop`. | | Scene-Entity | `Ferris wheel` is found in `Amusement park`. | |
| Entity-Concept | `BMW 320` is a kind of `Car`. | | Scene-Attribute | `Alleys` are `Narrow`. | |

TABLE 3: Examples of visual relations detected in NEIL [22]

*2) Statistic-based General Relation Extraction.* The statistic-based methods encode features such as visual, spatial, and statistics of the detected objects into distributed vectors and predict the relation between the given objects by a classification model. Unlike rule-based methods, statistic-based methods can detect all relations in the training set.

Some work has proved that predicting the predicates rely heavily on the categories of subjects and objects, but subjects and objects are not dependent on predicates, and there is also no dependency between subjects and objects [81]. For example, in triple (`Person`, `ride`, `Elephant`), `Person` and `Elephant` indicate that the relation might be `ride` rather than `wear`. Thus to utilize the dependency, [81], [82], [83] add language priors of language models into the statistic model by objects' labels and [84] set a stricter constraint that the hidden layer representation of a triple should satisfy *subject + predicate ≈ object*. It is embarrassing that the language model improves much, but the visual information contributes little [81].

Detected objects and relations in an image could be represented as a graph. The graph structure enables the edges to get more messages from other nodes and edges to classify the relation with higher accuracy. For example, [85] represents objects and relations as two complementary sub-graphs, where nodes are iteratively updated according to the values of the surrounding edges and vice versa. [86] used GCN to learn the context of objects and edges. Unfortunately, the recall@50 of triple detection in current visual detection models is still less than 23%, although the recall@50 of predicate detection has been up to 85.64% [87] on the visual relation detection benchmarks.

*3) Long-tail and Fine-grained Relation Extraction.* It is challenging for statistic-based methods to detect long-tail relations. Frequent relations are more likely to be predicted due to the bias of sample distribution in the training sets. Much work focuses on eliminating the effect of unbalanced samples in the training sets by metric learning [88], [89], transfer learning [90], few-shot learning [91] and contrastive learning [92], which are still limited to the feature fusion of hidden layers.

Fine-grained relation is a kind of long-tail relation. Existing studies on long-tail relation problems from the perspective of feature fusion fail to distinguish fine-grained relations well. For example, models tend to predict `on` instead of fine-grained relation `sit on/walk on/lay on`. For more informative unbiased predictions, [93] uses counterfactual causation instead of conventional likelihood to remove the effect of context bias. Differently, [94] orders relations in a hierarchy, from specific ones at the bottom to generic ones towards the top. It trains a classifier for each relation, classi-fying a detected triple into two types: whether it belongs to a certain relation or its sub-relations in the hierarchy.

**OPPORTUNITIES**. Despite much existing work, there still leaves many challenging issues unsolved. For instance: 1) *Visual Knowledge Relation Judgement*. Many visual triples extracted from images only describe the scene of the image, which are unqualified to be taken as visual knowledge since they are not widely accepted facts. The challenges (also opportunities) lie in how we recognize the triples of visual knowledge from the triples of scene information. 2) *Relation Detection based on Reasoning*. Existing relation detection methods predict the relations by a hidden unified representation fusing visual features and language priors. We cannot explicitly describe the basis of prediction. [95] builds a human action dataset to help predict an action by body part states. For example, if there is a person and a football in an image and (`Head`, `look at`, `Sth`) (`Arm`, `swing`, `-`) (`Foot`, `kick`, `Sth`) are meanwhile satisfied, the action will be judged as (`Person`, `kick`, `Football`). Unfortunately, this dataset is built manually. We need to summarize the chain of reasoning for relation detection automatically.

### 3.1.3  Visual Event Extraction

An event includes a trigger and several arguments with their argument roles. A trigger is a verb or a noun indicating the occurrence of an event. An argument role is a relation between an event and an argument, and the arguments are entity mentions, concepts or attribute values. The visual event extraction can also be divided into two sub-tasks: 1) to predict the visual event types; and 2) to locate and extract objects in source images or videos as visual arguments [23], [68], [96], [97]. This task is different from the situation recognition task [98], [99], [100] in CV, which aims to recognize a visual event rather than locating and extracting its visual arguments. Schemas defined in datasets of situation recognition tasks, such as SituNet [99] and SWiG [100], can be used to train models in this task.

**CHALLENGES**. The task has several challenges: 1) Visual event extraction requires pre-defined schemas for different event types, but there are a large number of visual events that experts have not defined. How to mine visual patterns as event schemas automatically? 2) How to extract visual arguments of a visual event from images or videos?

**PROGRESSES**. The existing work on visual event extraction mainly focuses on two aspects: 1) visual event schema mining, which detects and labels the most relevant visual entities (or concepts) as a new schema; 2) visual event

arguments extraction, which extracts argument role regions from visual data according to the event schema.

1) *Visual Event Schema Mining*. In large-scale visual event extraction, such as news, the visual schemas of many events have not yet been manually defined, which requires much experts' work. Large numbers of image-caption pairs from the web make it possible to mine and label the visual pattern for event schemas. Thus this task is reduced to finding a frequent itemset of visual patterns which indicate the correct event type from the images of a given event. The collection of images of an event can be retrieved from the image-caption pairs with the event's triggers as queries. Words or phrases in captions label the candidate image patches through visual grounding. Heuristic approaches (e.g., the Apriori algorithm) can be utilized to mine frequent visual image patches to find association rules for predicting the event type by visual patterns [96], [101].

Mining and labeling methods can correct wrong arguments or add missing ones in manually defined visual event schemas. For example, an ontology expert may consider `Explosion` and `Weapon` as important items in the schema of event `Attack`, but in some news corpus, these concepts are not discovered and `Smoke` and `Police` appears much more frequently, which is not expected in advance [101].

2) *Visual Event Arguments Extraction*. This task aims to extract a group of visual objects with the constraint of relations. The event types are classified according to the global features of images, and the event arguments are extracted as the most sensitive region to the event type by object recognition or visual grounding. The quality of the two sub-tasks on a large corpus is acceptable. In MMEKG [24], the instance-level evaluation has a precision score of about 64% on visual events and cross-modal triples.

In addition, the relations in visual and text arguments should also be aligned to ensure that the relations among visual objects are consistent with the relations in text. [68] aligns the situation graph [99] extracted from the image and the abstract meaning representation graph (AMR graph) [102] extracted from the caption of an event in terms of the semantics and categories of cross-modal arguments. Many constraints on semantic, event type, event argument role and the consistency between modalities are also added into joint extraction [23], [68].

Videos are more suitable for event extraction than images because the temporal bounding box of an event may be across the video, and all arguments may not appear in a single frame. [97] simplifies this task and extracts arguments from three keyframes derived from short video segments including only one event, and the keyframes are the most matching ones to the captions of the videos.

**OPPORTUNITIES**. The research on this task is still in an early stage, and many problems are still worth exploring. For instance: 1) The extraction of sequential events from a long video containing multiple events has not yet been addressed. 2) *Video Event Extraction with multiple Sub-events.* For example, the event `Making Coffee` is divided into a sequence of steps, such as `Cleaning coffee machine` → `Pour in the coffee beans` → `Turn on the coffee machine` and each step can be also considered as an event. The sequential steps need to be extracted

and listed by the timeline of the steps, which are difficult to be solved by current methods.

## 3.2 From Symbols to Images: Symbol Grounding

Symbol grounding refers to the process of finding proper multi-modal data items such as images to describe a symbol knowledge in a given KG, such as an entity, a concept or a relational triple. Some popular MMKGs constructed in the symbol grounding way are listed in Table 2(b).

In the rest of this subsection, we cover the process of grounding symbols to images in several fractionized tasks: *Entity Grounding* (Sec. 3.2.1), *Concept Grounding* (Sec. 3.2.2) and *Relation Grounding* (Sec. 3.2.3).

### 3.2.1 Entity Grounding

Entity grounding aims to ground entities in KGs to their corresponding multi-modal data such as images, videos and audios [12]. The existing work mainly focuses on grounding entities to their corresponding images.

**CHALLENGES**. The main challenges of grounding entities to images are the following: 1) How to find enough images with high quality for entities at a low cost? 2) How to select the images that best match an entity from much noise?

**PROGRESSES**. There are two major sources to find images for entities: (1) from *online encyclopedia* (such as Wikipedia), or (2) from the Internet through *Web search engines*.

1) *From Online Encyclopedia*. In Wikipedia, an article usually describes an entity with images. Wikipedia and DBpedia provide many facilities (such as Wikimedia Commons) to help build the connection between an entity in DBpedia and corresponding images or data in other modalities in Wikipedia. It is easy for researchers to use an online encyclopedia like Wikipedia to build the first version of a large-scale MMKG.

However, the encyclopedia-based approach has several major disadvantages: 1) First, not all entities are attached to many high-quality images in an online encyclopedia. We investigate that the average number of images per entity in Wikipedia is only 0.83. Second, many images of entities in Wikipedia are only indirectly related to that entity but can not accurately represent that entity. For example, there are several images of *animals*, *buildings*, *plaques*, *carvings* in images of `Beijing Zoo` in Wikipedia. Third, the images of the non-visualizable entity may bring mistakes. For example, in the Wikipedia article of `Gaussian Progress`, there is an image of *Gaussian processes with different prior conditions*, which should not be mapped to any image. Finally, the coverage of MMKG built from Wikipedia alone still needs to be improved. English Wikipedia has 6 million entities (articles), which is the upper bound of the capacity of the MMKG harvested from English Wikipedia. According to our investigation, 79.35% of Wikipedia articles in English have no corresponding images, and only 6.7% of them have at least 3 images.

2) *From Search Engines*. Search engine based solutions are proposed to improve the coverage of an MMKG. We can easily find images from the search results of a commercial search engine by specifying entity names as queries, where the top-ranked image is more likely to be the correct image

| Natalie Portman | Keira Knightley | fire fighter | trash collector |
| (a) | | (b) | |

Fig. 6: Examples that can hardly be distinguished by visual entity extraction methods. (a) Similar visual entities: `Natalie Portman` and `Keira Knightley`; (b) Similar visual concepts: `fire fighter` and `trash collector`.

S1: In 1964, *Trump* enrolled at Fordham University.
S2: In 1971, *Trump* was named president of the family company and renamed it The Trump Organization.
S3: *Trump* registered as a Republican in Manhattan in 1987.
S4: *Trump* is the wealthiest president in U.S. history, even after adjusting for inflation
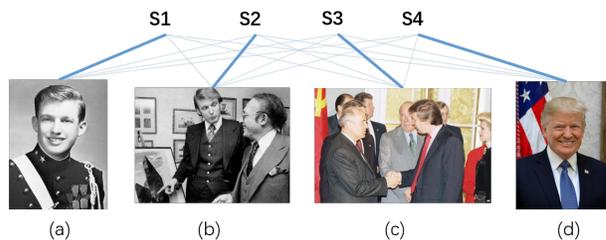


Fig. 7: Take `Trump` as an example to illustrate that an entity needs different images to express its different aspects (Trump as (a) a young student, (b) a businessman, (c) a politician, or (d) the president of the USA) in different contexts

of the searched entity. Thus we can select these images for the entity to be searched. Compared to the Wikipedia based approach, the coverage of MMKG is significantly improved in the search engine based approach.

However, the search engine based approach is easy to introduce noisy images into MMKGs. It is well recognized that the search engine results might be noisy. Another reason is that it is not trivial to specify the search keywords. For example, the search query *"Bank"* is not good enough to find the image for `Commercial Bank`, since it also incurs the images of `River Bank`. Hence, many efforts have been made to clean candidate images. The query words are usually extended for disambiguation by adding parent synsets [103] or entity types [28]. Diversity is also a non-negligible issue when selecting the best images for the entity. An image diversity retrieval model is trained to remove similar redundant images so that the grounded images are as diverse as possible [29].

Compared to the encyclopedias-based approaches, search engine based approaches are better in coverage but worse in quality. The two approaches are often used together since in most cases the knowledge acquired by these two approaches complements each other [29]. For example, the coverage of MMKG harvested from Wikipedia can be improved by collecting more images for each entity from search engines [29].

Due to the decoupling of entities and their visual features, an MMKG constructed with encyclopedias or search engines can distinguish visually similar entities, as shown in Fig. 6. Entity grounding methods make it possible to build a domain-oriented fine-grained MMKG (e.g., a movie/product/military MMKG).

**OPPORTUNITIES**. There are many unsolved problems in this direction. 1) Entities are grounded into several images, each of which is only an aspect of the entity. For example, the image collection of a person may be images of different ages, life photos, event photos, single photos and family photos. How do we determine the most typical subset?2) Real-world entities are multi-faceted, and it is desirable to associate an entity with multiple images in different contexts. The demand motivates us to propose a new task *multiple grounding* that selects the most related images from the entity given a specific context. For example, `Donald Trump` has a lot of different images that can be collected from the web. But as shown in Figure 7, any single image is not appropriate for all the different contexts. Thus, `Trump` should be multi-grounded when constructing the

knowledge graph. 3) If there is an objective domain corpus containing a large amount of texts with attached images, we may convert the entity grounding task into a text-image retrieval task, such as the work done on the E-commerce domain [20].

### 3.2.2 Concept Grounding

Concept grounding aims to find representative, discriminative and diverse images for visual concepts.

**CHALLENGES** Although some visually unified concepts (such as `man`, `woman`, `truck` and `dog`) can also be grounded to images with the entity grounding methods introduced in Sec. 3.2.1, the symbol grounding to the other concepts faces new challenges: 1) Not all the concepts could be adequately visualized. For example, `irreligionist` cannot be grounded to a specific image. How to distinguish visualizable concepts from non-visualizable ones? 2) How to find representative images for a visualizable concept from a group of relevant images? Note that the images of a visualizable concept might be very diverse. For example, when it comes to `Princess`, people often think of several diverse images: *Disney princesses, ancient princesses in historical movies or modern princesses in the news*. Therefore, we have to consider the diversity of images..

**PROGRESSES**. In response to the above challenges, related studies are divided into three tasks: visualization concept judgment, representative image selection and image diversification.

1) *Visualization Concept Judgment*. The task aims to automatically judge visualizable concepts and is a new task to be solved. [104] discovers that only 12.8% of the synsets of `Person` subtree have well-accepted imageability (i.e., the score is greater or equal to 4 and the total score is 5), and many of the rest synsets have no corresponding visual descriptions. For example, `Rock star` is imageable, and `Job candidate` is non-imageable. So what are the criteria for recognizing visual concepts? The manual annotation in [104] is unpractical in constructing a large-scale MMKG.

In order to automatically judge visual concepts, there has been much effort based on syntax and semantics. [105] thinks that abstract nouns concepts are non-visualizable

so that TinyImage dataset [105] removes all hyponyms in the subtree of `Abstraction` in WordNet and only collects images for non-abstract noun concepts. However, these methods are not very accurate. For example, `Anger` or `Happiness` can be grounded in an image of a person who feels angry or happy. Since the images come from the web, it is possible to use search engine hits to judge visual concepts. For example, a word might be visualizable if the number of Google image hits is larger than that of Google web hits [106]. [107] assumes that if images of a concept from Google are similar (with a small variance), this concept is more likely to be visualizable. This assumption may lead to a low recall, so it is used to correct the false negative predictions (non-visualizable) of classifiers.

2) *Representative Image Selection*. Based on the methods of Sec. 3.2.1, we get a collection of images for each visual concept. This section focuses on selecting visually representative and discriminative images in the collection.

The task aims to re-rank the images according to their representativeness. The representative scores of images derive from results of cluster-based methods, such as K-means, spectral clustering, etc. The smaller the variance within a cluster, the higher the scores of images in the cluster. After re-ranking the representative scores of images, the top may be representative images. In addition, the expected images are also constrained by rules to distinguish different clusters. For example, [108] adds a new metric to rank images together with similarity within clusters, which is the ratio of inter-class distances and intra-class distances, and the bigger a ratio, the more discriminative the image is.

The captions and tags of images from search engines could also be utilized to evaluate the representativeness and discrimination of images at the level of semantics. Captions and tags provide semantic information that images do not have. For example, a photo of *Icelandic landscapes* and a photo of *British landscapes* may look similar, but text tags can help us distinguish their differences in concepts. In [106], [109], [110], tags are clustered based on semantic features and images are reassigned into each cluster according to their tags' semantic clusters.

3) *Image Diversification*. The task requires that images in which concepts are grounded should balance diversity and relevance. The images should also be re-ranked after clustering, but the difference from representative image selection is that we want to show the results of as many clusters as possible. Specifically, in each selection step, images from unselected clusters are preferred to be selected.

There are two types of scores for ranking the priority of selection: diversity scores and relevance scores, where diversity scores evaluate the topics of images and relevance scores penalize the difference of images to avoid semantic drift. For fusing the two conflicting scores, [111], [112] use Max-Min methods to choose candidates: assign a higher score to images that are not similar to the selected set, and choose the dissimilar one with the highest score among the remaining similar ones. [113] mines topics (e.g., `View`, `Flag`, `Map`) from image captions of popular entities (e.g., `Greenland`) to expand queries of long-tail entities of the same type (e.g., `Country`) during image retrieval. Then images of long-tail entities are filtered by local outlier factors based on the distribution of similar popular entities' images.
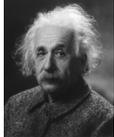
| concept type | visualizable concept | non-visualizable concept |
|---|---|---|
| example | Surgeon | Physicist |
| image | | |

TABLE 4: Examples of visualizable concept grounding and non-visualizable concept grounding. The visualizable concept `Surgeon` can be grounded to *the photo of doctors wearing surgical suits and performing surgery in the operating room*, and the non-visualizable concept `Physicist` can be grounded to *the photo of Einstein* since `Einstein` is a typical entity of `Physicist`.

Diversity is achieved by pattern mining, and relevance is achieved by pattern transferring.

We can also resolve the ranking problem by graph algorithms. A set of images could be represented as a graph, where images are nodes and visual similarities between images are weights of edges. Thus, the ranking of representative images reduces to finding an optimal path in a fully connected graph concerning re-weighted values of edges. [114] uses dynamic programming to search for the optimal sequence in an image graph, where the value of edges is a joint criterion combining diversity score and relevance score. Markov random walk is also used for the optimal sequence in [106], [115], where [115] weights the values by Max-Min methods and [106] reassigns the visits values between nodes according to their source clusters by a two-layer graph model.

These studies concentrate on text-image retrieval, and only [113] is related to MMKGs. There are still many unsolved biases on the diversity of images of concepts derived from the Internet on gender, race, color and age, and the problem now relies heavily on crowdsourcing [104].

**OPPORTUNITIES**. As a fledgling area, many unsolved problems are left for future research. We give two examples below:

1) *Abstract Concept Grounding*. Previous work on concept visualization judgment seldom considers abstract concepts. But the abstract concepts could also be grounded in images. For example, `Happiness` are usually associated with *smile*, and `Anger` are usually associated with *an angry face*. Some abstract nouns have a diverse but fixed visual association, such as nature, human and action. For example, in [116] the images of `Beauty` are associated with following word clusters: *woman/girl*, *water/beach/ocean*, *flower/rose*, *sky/cloud/sunset*. Similarly, the image of `Love` are associated with following word clusters: *baby/cute/newborn*, *dog/pet*, *heart/red/valentine*, *beach/sea/couple*, *sky/cloud/sunset*, *flower/rose*. It shows that some abstract nouns often have generic and fixed images in terms of sentiment and discriminative images in terms of semantics.

2) *Gerunds Concept Grounding*. Gerunds are a special kind of nouns that could be transformed into verbs, such as *singing* → *sing*. [80] grounds many gerunds to images through crowdsourcing, such as `arguing with`, `wrestling with` and `dancing with`. These verbs about

human interaction are sensitive to the features of body angle, gaze angle, the position of the joints and expression.

3) *Non-visualizable Concept Grounding via Entity Grounding*. If a concept is non-visualizable but its hyponym entities could be visualized, the concept could also be grounded via its entities. For instance, a reasonable selection of the grounded image for such a concept is to use the image of the concept's most typical entity. As shown in Table 4, we use a photo of `Einstein` to ground the concept `Physicist`. It is reasonable since most of us will think up with `Einstein` when we mention a `Physicist`. However, there are still a lot of unresolved questions: (a) In general, different people will come up with different typical entities for a concept, so we should address such subjectivity in concept grounding. Whether an entity is a typical one in the constrain of its concept? (b) We should choose several typical entities' images to present that concept. How do we summarize and select typical entities to represent concepts? (c) Whether should we abstract common visual features from multiple images of entities?

### 3.2.3 Relation Grounding

Relation grounding is to find images from an image data corpus or the Internet that could represent a particular relation. The input could be one or more triples of this relation, and the output is expected to be the top-ranked representative images for the relation. For example, (`Justin Bieber, couple, Selena Gomez`) could be grounded to an image of "*Selena Gomez and Justin Bieber Kissed*" instead of "*Selena Gomez and Justin Bieber worked out together*".

**CHALLENGES**. When we take a triple as a query to retrieve images for the relation, the top-ranked images are often more relevant to the subject and object of the triple but not to the relation itself. How to find images that could reflect the semantic relation of the input triples?

**PROGRESSES**. Existing efforts on relation grounding mainly focus on the co-occurrence of visual objects in images or textual entities in captions. Richpedida [29] proposes a very strong assumption that if there is a pre-defined relations (e.g., `nearBy` and `contain`) between two entities in the Wikipedia descriptions, the same relations also exist between two entities' corresponding visual objects. But in reality, it is more likely that the two objects do not simultaneously appear in one image. Even if they do, the relation shown in the image may not be the expected one.

Relation grounding could be modeled as a fine-grained text-image retrieval problem, where the triple *(subject, relation, object)* is the query and candidate images are represented with the implicit or explicit structure information of the scene graph extracted. Specifically, each image could also be represented as a combination of multiple $(s, p, o)$ by multi-branch CNN [117] or graph convolutional neural network (GCN) [118].

Instead of global matching of cross-modal embeddings, we expect the item-by-item matching of objects and relations. If we represent the textual query and candidate images into graphs, the relation grounding task turns into a task of graph matching, as illustrated in Figure 8. [119] represents the two graphs by GCN, in which objects are updated from themselves and relation nodes are updated
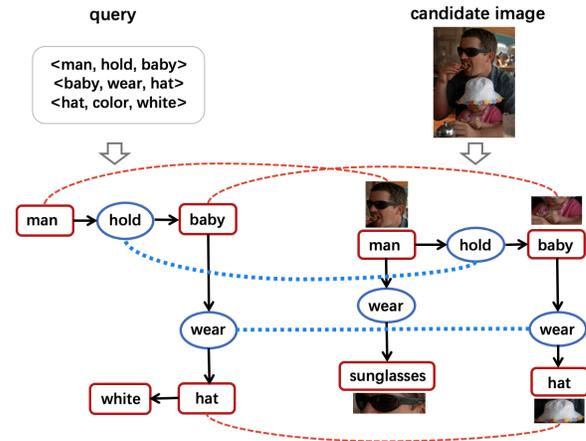


Fig. 8: Relation grounding is often considered as a fine-grained text-image retrieval problem. The queries are one or more triples, and the expected images should be consistent with entities and relations in the query. The figure shows an example of relation grounding by graph matching in [119].

from the aggregations of their neighbors. In predicting, the similarity between two graphs is measured by matching object nodes and relation nodes, respectively.

**OPPORTUNITIES**. Existing studies mainly focus on grounding spatial relations and action relations, such as `leftOf`, `on`, `ride` and `eat`, which could be observed visually in images. However, most semantic relations such as `isA`, `Occupation`, `Team` and `Spouse` may not be that visually obvious in images. There is a lack of training data for these relations, thus it is difficult to train models to retrieve images with the above solutions. Fortunately, some datasets [120], [121] of relation extraction based on textual named entities and visual relations could be helpful.

## 3.3 Comparing Two Construction Ways

There are several differences between the image labeling and symbol grounding solutions for constructing an MMKG in the aspects of applicable scenarios, construction efficiency, quality, etc. We analyze kinds of MMKGs in which multi-modal data are not only images but also code, audio or video, and summarize these differences as follows:

1) *Applicable Scenarios*. If the multi-modal data are treated as first-class citizens in some scenarios, the multi-modal data labeling way is more preferred to construct the MMKG, such as unearthed oracle bones' photos in oracle bones recognition system [136], teachers' class audios in educational services [137] and the movies' videos in deep video understanding tasks [94]. If the multi-modal data collected is redundant and noisy, the multi-modal data labeling way may produce many low-quality (such as repeated or mismatching) visual entities. In this case, the symbol grounding way is preferred to construct the MMKG because the symbols in KGs have already been well filtered and refined, such as the movie ontologies in recommendation systems [19], product ontologies in e-commerce dialogue systems [20] and paper ontologies in academic information retrieval and KBQA [138], [139].

| Multimodal Application | Benchmark Datasets | Advantages with MMKG |
|---|---|---|
| Entity Recognition and Linking | Twitter2015 [50] Twitter2017 [122] Weibo [123] WikiDiverse [124] | 1.background knowledge provides deep features of images 2.images provide necessary complementary information, help to capture the relationship among mentions and entities 3.learn distributed representations for each entity with multi-modal data |
| VQA | GQA [125] OK-VQA [49] FVQA [126] KVQA [127] KB-VQA [128] | 1.provide knowledge about the named entities and their relations in the image, leading to deeper visual content understanding 2.conduct the reasoning process and predict the final answers in a more explicit way with symbolic knowledge from MMKG 3.refine the answers with more interpretability and generality |
| Image-text Matching | Flickr30k [59] MSCOCO [58] Visual Genome [53] | 1.expand more semantic concepts 2.introduce informative relationships between visual concepts by constructing scene graphs 3.enhance the reasoning capabilities of multi-modal data with graph-structured information |
| Image Tagging | NUS-WIDE [129] | help disambiguation the concept and relate them better to images |
| Image Captioning | MSVD [130] MSCOCO [58] GoodNews [131] | 1.enable the understanding of unseen objects with MMKG symbolic knowledge 2.leverage MMKG for relational reasoning to generate more accurate and reasonable captions 3.capture fine-grained relationships between entities in different modalities |
| Visual Storytelling | VIST Dataset [132] | 1.triples in MMKG provide explanation and traceability for described facts 2.provide a strong logical inference between images for more fluent story |
| Recommender System | MovieLens [133] IntentBooks [134] Dianping [19] KKBOX [135] | 1.provide background knowledge for items with rich semantics to solve the cold-start problem 2.learn through rich path semantics across different modalities in MMKG and produce an interpretable and explicit recommendation 3.construct personalized MMKG for items and model entity relation reasoning between them |

TABLE 5: Benchmark datasets for their corresponding multimodal applications incorporating MMKGs.

Whether multi-modal data or symbolic knowledge is first-class citizen depends on what kind of knowledge we want the MMKG to provide. For example, in [139] when we want to know the relations between geoscience academic papers and maps in them, the papers are first-class citizens; when we want to know the relations between maps and regions pointed in these maps, the maps are first-class citizens.

2) *Efficiency.* The symbol grounding solutions are usually retrieval-based methods [20], [28], [29], [30], [140], [141] and the multi-modal data labeling solutions are usually classification and detection methods [21], [23], [24], [94], [136]. Extracting entities, concepts and relations in multi-modal data labeling solutions is time-consuming [22]. Therefore, it will be an excellent choice to start the construction of an MMKG from scratch using the symbol grounding solutions.For example, NEIL [22] initially collects image datasets by retrieving images from search engines with ontologies of NELL [142] as queries and then extracts objects and relations in these images.

3) *Quality.* Except for the quality of extraction models, the multi-modal data labeling solutions have to solve the problem of coarse-grained labeling and inappropriate semantic hierarchies. Symbol grounding solutions could solve these problems. However, the symbol grounding way also faces the problem of missing and mismatching images of symbols. For example, it is easy to find a bad image for a long-tail entity from search engines. Because such an entity might have no image on the web, any clicked image is misleading to a mistake grounding.

# 4 APPLICATION

After a systematic review of MMKG construction, this section explores how the knowledge in MMKGs can be applied to and benefit a wide variety of downstream tasks. For a quick overview, Table 5 lists some mainstream application tasks, their benchmark datasets, and the advantages brought by MMKGs. We categorize such tasks into (i) in-KG applications (Sec. 4.1) , (ii) out-of-KG applications (Sec. 4.2) and (iii) domain applications(Sec. 4.3), discussed as follows.

## 4.1 In-MMKG Applications

In-MMKG applications refer to tasks conducted within the scope of the MMKG where the embeddings of entities, concepts and relations are already learned. Thus, before introducing in-MMKG applications, we briefly go through the distributed representation learning of the knowledge in MMKGs, also named MMKG embedding.

The MMKG embedding models are developed from the embedding models on conventional KGs, i.e., *semantic matching based* models, RESCAL [143] and its variants [144], [145], which measure the possibility of existence of triple ($h$, $r$, $t$) by the calculation of $h$, $r$, $t$ in vector space, and *translational distance based* models, TransE [146] and its variants [147], [148], [149], which should conform to the assumption: $t \approx h + r$. $h$, $t$, $r$ is respectively the vector representation of head entity, tail entity and relation in a triple. There are two additional issues in dealing with multi-modality data: how we effectively encode the vision knowledge and information contained in images, and how we fusion knowledge of different modalities. 1) *Vision Encoders.* With the development of deep learning, hidden features gotten from CNN [145], [150], [151] or Transformers [152] are the main image embeddings used in MMKG representation, while other explicit visual features such as GHD, HOG, CLD can hardly be leveraged in MMKG representation. 2) *Knowledge Fusion.* There are two ways to fuse the knowledge embeddings of multi-modalities: combining every single modal representation trained in its own vector space (such as concatenation, average pooling, SVD and PCA) [27], [28], [151], or further learning a unified embedding by projecting different modal representations into the same space [145], [150], [153]. While some methods [151] take the fused results as the MMKG embedding directly, the other methods [145] further train the uni-modal representations on a well-designed objective function.

In the following, we introduce four well-studied in-MMKG applications including *link prediction*(Sec. 4.1.1), *triple classification*(Sec. 4.1.2), *entity classification*(Sec. 4.1.3), and *entity alignment*(Sec. 4.1.4).

### 4.1.1   Link Prediction

Link prediction in MMKG [150], [153] aims to complete a triple $(h, r, t)$ when one of the entities in $h, r, t$ is missing, i.e., predicting $h$ in $(?, r, t)$ or predicting $t$ in $(h, r, ?)$. A similar task is to predict the missing relation between two given entities, i.e, predicting $r$ in $(h, ?, t)$.

Conventionally, link prediction on KGs can be processed with a simple ranking procedure, which finds the best fit entity to complete a triple from all the candidate entities. Specifically, in the training stage, the embedding model learns an embedding for each entity or relation, for instance, with the training objective $\boldsymbol{t} \approx \boldsymbol{h} + \boldsymbol{r}$ defined by TransE [146]. Then in the prediction stage, the most matching $h$ in $(?, r, t)$ is found by ranking all candidate head entities $h^*$ according to a score function like $\arg\max_{h^*} \phi(h^*, r, t)$, where the score function is diverse in different embedding models [154].

Compared to the task in traditional KGs, the images fused into representations of entities and relations in MMKGs could provide extra visual knowledge to enrich the information of embedding. For instance, the images of a person might provide evidence for the person's age, profession, and designation [145].

This task is different in existing MMKGs depending on the scenario. IMAGEgraph [27] proposes to express the relation prediction between unseen images and multi-relational image retrieval as visual-relational queries, such that these queries could be leveraged for MMKG completion. Compared to the conventional way, IMAGEgraph performs better on the relation and head/tail entity prediction tasks and is able to be generalized to unseen images, to answer some zero-shot visual-relational queries. For example, given an image of an entirely new entity not part of the KG, this approach can determine its relation with another given image for which we do not know the underlying KG entity.

Similarly, MMKG [28] constructs three datasets to predict the multi-relational links between entities, with all the entities associated with numerical and visual data. However, it only focuses on the `sameAs` link prediction task and answers such queries for MMKG completion. Three quite heterogeneous knowledge makes MMKG a vital benchmark to measure the performance of multi-relational link prediction methods and validates the hypothesis that different modalities are complementary for the `sameAs` link prediction task.

### 4.1.2   Triple Classification

Triple classification aims to distinguish correct triples from incorrect ones, which can also be seen as a sort of KG completion task. Based on the embedding model learned on an MMKG, each triple could be calculated with an energy score $E(h, r, t)$. Different threshold $\delta_r$ is set for each relation $r$, and a triple will be predicted to be negative if its energy score is higher than $\delta_r$. In classification models, correct triples are corrupted by replacing one of the $h, r, t$ to generate negative data [150], [153].

### 4.1.3   Entity Classification

Entity classification categorizes entities into semantic categories, i.e., concepts of different grains in the MMKG. Entity classification can also be regarded as a special link prediction task, where the relation is `IsA` and the tail of the triple to be predicted is a concept in the MMKG.

Various entity classification models have been proposed for traditional KGs, which could also be adopted in MMKGs. But the rich multi-modal data for entities and concepts in MMKGs cannot be fully utilized without a good MMKG embedding model. For instance, some efforts [140], [155] work on learning embeddings for entities and concepts from several different types of modalities and then encode them to a joint representation space. However, [140] argues that this task in KGs cannot be solved purely by node embedding models, and the graph structures should also be considered. Therefore, [140] proposes a collection of extensive and high-qualified multi-modal benchmarks for precisely evaluating node classification tasks on MMKGs.

### 4.1.4   Entity Alignment

Entity alignment works on aligning entities that refer to the same real-world identity in different MMKGs. It is a viable way to integrate two MMKGs into one when there are overlaps.

The core idea is to learn representations for entities in different KGs and then evaluate the similarity between each entity pair between the two KGs. The features used in entity embedding between two traditional KGs include in-KG context information (e.g., the semantics of OWL properties, co-occurrence of neighbors, compatible attribute values) and external information (e.g., external lexicons and Wikipedia links). For MMKGs, due to the introduction of multi-modal features, some entity-alignment oriented MMKG embedding models are proposed [156], [157]. Feature vectors are encoded for different modalities respectively and then merged into one to represent the entity by the knowledge fusion techniques mentioned at the beginning of this subsection. One work [156] uses ranking loss as the loss function, while another [157] designs a loss function $L = \alpha||e - e_s|| + \beta||e - e_n|| + \gamma||e - e_i||$ to enhance the complementarity of multiple modalities, where $e_s, e_n, e_i$ is the embedding of three different modalities respectively $e$ is the final embedding of the entity, and $\alpha, \beta, \gamma$ is ratio hyper-parameters for each modality.

Another line of work [28] elaborates a Product of Experts (PoE) model to answer queries such as $(h?, sameAs, t)$ or $(h, sameAs, t?)$ where $h$ and $t$ are from different KGs. By incorporating [158] and extending it to visual features, the end-to-end learning framework is superior to the concatenation and an ensemble type of approach for entity alignment.

## 4.2   Out-of-MMKG Applications

The out-of-KG applications refer to the downstream applications that are not limited to the boundary of MMKGs but could be assisted by them. In the following, we introduce several such applications as examples. Instead of providing a systematic reviews to all the solutions of these tasks, we mainly focus on introducing how MMKGs are utilized, and the advantages of MMKGs compared with other solutions.

### 4.2.1   Multi-modal Entity Recognition and Linking

Named entity recognition (NER) with plain texts has been studied extensively. Ambiguity and diversity of entity men-

tions have always been the key challenges. Recent work focusing on detecting entities from texts attached with images is defined as multi-modal NER (MNER) [50], [122], where images could provide necessary complementary information for entity recognition.

MMKGs can enhance MNER by providing vision features of entities to enhance the representation of images or text. For instance, [159] compares the given image with the images of candidate entities (from text) and two-hop neighborhood entities in the MMKG to find the most relevant entity as external background knowledge for disambiguation. [124] also employs MMKGs to retrieve more labels as related words based on the co-occurrence frequency between entities. With the expansion of entity type labels from MMKGs, more task-specific salient features are highlighted, avoiding being neglected in cross-modal interactions and improving the performance of MNER.

Given a text with images attached, multi-modal entity linking (MEL) uses textual and visual information to map an ambiguous mention in the text to an entity in a given KG [160]. Although some early efforts do MEL based on a traditional KG, increasingly recent work uses MMKGs for linking. MEL utilizes the knowledge with images in an MMKG in two ways: (1) providing the target entities to which the entity mentions should be linked; (2) learning distributed representations for each entity with multi-modal data, which are then used to measure the correlation between a mention and an entity. The usage of visual information with images would help to capture the relationship among mentions and entities [160], [161], but the irrelevant part with images may also become noises and bring negative impact to the representation learning for both mentions and entities. To remove the side effect, a two-stage image and text correlation mechanism is proposed to filter out the irrelevant images based on the pre-defined threshold, and the multiple attention mechanisms are also utilized to capture the critical information in the mention representation and entity representation by querying multi-hop entities around the mention's candidate entities [123].

### 4.2.2 Visual Question Answering

Visual question answering (VQA) is challenging, requiring accurate semantic parsing of the questions and an in-depth understanding of the correlations between different objects and scenes in the given image. In most recent VQA benchmark datasets such as GQA [125], OK-VQA [49] and KVQA [127], many questions require visual reasoning combined with external knowledge. The newly proposed VQA tasks bridge the discrepancy that humans can easily combine knowledge from various modalities to answer visual queries. For example, in the question *"Which American President is associated with the stuffed animal seen here?"*, if the stuffed animal in the image is detected as "`Teddy Bear`", the answer inferred through KG will be "`Theodore Roosevelt`", who is often referred as *"Teddy Roosevelt"*, and after whom *Teddy Bear* is named [49].

Obviously, reasoning only by semantic parsing and matching can not answer the above question [128]. In this case, MMKGs could help in three aspects. First, MMKGs provide external knowledge about the named entities and their relations in the image, leading to deeper visual content understanding. Second, the facts about visual entities in the image and textual entities in the question from existing MMKGs help to re-weight the answer [162], which also benefits from the unified representation of all modal resources including images, questions and structured facts. Third, entities and relation triples of different modal in MMKGs can be represented as nodes and edges in a heterogeneous graph and represented in a unified format, which facilitates explicit reasoning with heuristic rules, SPARQL queries [128] or weighted passing messages between GNN nodes [51], [128], [162].

Some recent efforts tend to construct MMKGs for VQA by combining existing KGs and well-annotated image datasets. For example, the explicit knowledge in [51] has four sources: `hasPart` triples from hasPart KB [163], `hasPart/isA` triples from DBpedia [6], commonsense triples from ConceptNet [2], and location triples of visual objects from Visual Genome [53]. The model fusing explicit symbolic knowledge from the MMKG and implicit knowledge from VL-PTMs outperforms the pure VL-PTMs, and most of the knowledge in the MMKG is non-overlapping with the implicit knowledge in VL-PTMs [51].

### 4.2.3 Image-Text Matching

Image-text matching is a fundamental task in many cross-modal applications like image-text and text-image retrieval, which aims to output a semantic similarity score between the input image and text pair [164], [165], [166], [167], [168].

Image-text matching is usually achieved via mapping texts and images into a joint semantic space and then learning unified multi-modal representations for the similarity calculation. A general method is to exploit a multi-label detection module to extract semantic concepts and then fuse these concepts with the global context of image [165], [169], [170]. However, it is difficult for pre-trained detected-based models to find long-tail concepts, which constrains models to those detected concepts and leads to poor performance.

To overcome the bias in the training data for retrieval tasks, MMKG could be leveraged to expand more visual and semantic concepts leveraging the relations between multi-modal entities. Besides, MMKGs can also help to construct scene graphs, which introduce informative correlation knowledge between visual concepts and further enhance image representations. For example, the concept pairs that frequently co-occurred in the multimodal triples of an MMKG, such as `house-window` and `tree-leaf`, can be extracted to enhance the representation of concepts in images, thus providing a solid context signal for semantic understanding of images and leads to improved performance of image-text matching [16]. Besides, considering that one key step in the image-text matching task is to align both local and global representations across different modalities, some efforts propose incorporating relations in MMKGs to represent both image and text with higher-level semantics [171]. Such graph-structured information better enhances the reasoning and inference capabilities of multi-modal data with more interpretability. MMKG also helps cross-modal alignment by learning a more unified multimodal representation.

### 4.2.4 Multi-modal Generation Tasks

Several vision-text generation tasks, such as image tagging, image captioning, visual storytelling, etc., could benefit from MMKGs.

**Image Tagging**. Traditional image tagging methods are limited by biased distribution, noise and imprecise tags. MMKGs not only establish a well-organized taxonomy of concepts (such as synonyms, hypernyms and hyponyms) but also provide corresponding representative and discriminative images for concepts, thus they could greatly alleviate the effects of distribution bias of tags and noisy tags. For example, [172] constructs an MMKG called VTKB containing hierarchical concepts, linking concepts of original tags to images and linking images by the similarities of embeddings. The candidate concept set is a subset of the union of the parent, the child, the part, the whole, synonyms, hypernyms, hyponyms and related concept sets of the original coarse-grained tags of images. Finally, the re-generated fine-grained tags are those concepts that best match nearest neighbor images, where the candidate concept set depends on the type of bias specified in advance. The experimental results show that the proposed method with MMKGs achieves higher mean average precision than the baselines without MMKGs. MMKGs help to generate more relevant candidate tags and are more capable of disambiguating them than ConceptNet, WebChild and ImageNet.

**Image Captioning**. The mainstream statistic-based image captioning models have two weaknesses: First, they heavily rely on the performance of object detectors. The encoder-decoder framework with separate procedures of detection and captioning always leads to semantic inconsistency between the pre-defined objects/relations and target textual descriptions. Second, unseen objects always pose great challenges to them. The models trained on image-caption parallel corpora always fail to describe unseen objects and concepts.

Fortunately, MMKGs could help to alleviate the two obstacles in the following ways: 1) Some efforts [173] propose to leverage MMKG for relational reasoning, which results in more accurate and reasonable captions. More specifically, a semantic graph could be built for visual and knowledge vectors embedded from candidate image proposals, and the semantic graph could then be encoded for textual description generation. In this way, the semantic constraints summarized in MMKGs can be fully used, which may further endow the MMKGs ability and readily extended for more advanced reasoning. 2) The symbolic knowledge from MMKGs may enable the understanding of unseen objects [48], which are made visible by the semantic relation between seen objects and unseen objects in MMKGs. In the knowledge-guided image-caption task containing novel objects, the key module is a multi-label image classifier for grounding depicted visual objects to knowledge base entities, unveiling a way to build a connection between real-world objects to their multi-modal information with the assistance of MMKGs [48]. By introducing external knowledge from an MMKG-based multi-label classifier, image representations are also expanded.

A more complex task, named entity-aware image captioning, asks for more informative descriptions of named entities based on the background knowledge in the given article. In this task, these methods that only focus on textual knowledge and neglect the associations between named entities and visual cues in the image perform badly. However, MMKGs are very handy for the task requiring fine-grained cross-modal alignment between named entities and their images and further extension. In [18] the textual scene graph and visual scene graph extracted from the input article and images are aligned by the cross-modal entity matching module pre-trained on Wikipedia articles and images. Incorporating the aligned cross-modal scene graphs and external knowledge from Wikipedia, more accurate named entities and relevant events are chosen and refined. The results show that the structurization of cross-modal data improves the value of BLEU, METEOR, ROUGE, CIDEr and entity F1, where structurization with external knowledge significantly improves the performance.

**Visual Storytelling**. Visual storytelling is more challenging, aiming to tell the story according to several successive images. This task requires discovering the relations between the images and the objects associated with the images. Traditional visual storytelling approaches usually treat the task as a sequential image captioning problem and ignore the relation between images, which may produce monotonous stories. Besides, these approaches are limited to the vocabulary and knowledge in a single training dataset. To tackle these problems, [174] resorts to an MMKG for help within a distill-enrich-generate three-stage framework. After extracting a set of words from each image, all words from two consecutive images are paired to query the MMKG (such as Visual Genome) to enrich possible triples. Then story sentences are generated based on the most reasonable triple step by step. The methods using the relations in KGs show a strong ability of logical inference between images, generating more fluent stories than non-KG methods, and the triples from Visual Genome perform better than those from OpenIE in this task.

### 4.2.5 Multi-modal Recommender System

Recommender systems aim to recommend items that users might like/buy through the analysis of historical data, where accuracy, novelty, dispersity, stability and other factors should be balanced [175], [176]. Where there are multi-modal data such as image and text in a recommending scenario, we say it is a multi-modal recommender system, where the information of different modalities should be leveraged jointly.

It has been proved that MMKGs could greatly enhance multi-modal recommender system [177]. First, MMKGs incorporate different modal data with a hierarchical structure, enriching the representations of items [19], which can be used to solve the cold-start problem long existing in collaborative filtering based on recommending strategies [178]. Second, MMKGs can be used to select better logical reasoning paths for more explicit and explainable recommendations. For instance, [179] takes advantage of the the graph structure of MMKGs to design a hierarchy-based attention-path, which reduces the size of the action space and lets the model be more focused on critical intermediate items (entities). The results imply that additional structured textual and visual knowledge can significantly improve the recommendation quality [19], [178], [179].

| entity | Pluvianus aegyptius | The Wandering Earth | arrogance |
|--------|---------------------|---------------------|-----------|
| image  |  |  |  |

TABLE 6: Examples of quality problems in MMKG, such as images of frequently co-occurring entities, long-tailed entities, and abstract entities.

## 4.3 Domain Applications

In addition to applications on movie recommender [19] or e-commerce KBQA systems [20], MMKGs are also applied in multi-modal tasks such as cross-modal retrieval, dialogue system and object detection in some domain applications. For instance, [139] uses a geoscience academic MMKG to help to retrieve multi-hop queries, such as papers about specific geographic locations with a certain affiliation. [138] uses an academic MMKG about papers and codes to offer retrieval on the implementation level. In some other works, MMKGs are adopted to enrich the representation of entities with the help of images (e.g., X-rays, CT and ultrasound) and textual description, improving the performance of doctor-patient dialogue systems of COVID-19 [141] and further reducing the risk of close contact. In the archaeology field, MMKGs also contribute to oracle bones detection and recognition, not only taking into account edges, textures, cracks, scratches, splinters and background, but also offering relevant literature, location and institutions to assist decision making [136].

## 5 OPEN PROBLEMS

### 5.1 Complex Symbolic Knowledge Grounding

Besides entities, concepts and relations, some applications require the grounding of complex symbolic knowledge consisting of multiple relational facts with close semantic relations. These multiple relational facts may be a path or a subgraph in a KG. For example, for a subgraph in a KG containing Trump's wife, daughter, grandson etc., a proper grounding image might be a *Trump's family* photo. This motivates *multiple relational grounding*, which aims to find images to express the knowledge in a path or a subgraph in a KG. Multiple relational grounding is challenging since it involves the grounding of more than one relation, which is usually interleaved with each other in a complicated way.

### 5.2 Quality Control

Besides the common quality problems studied extensively in traditional KGs (e.g., accuracy, completeness, consistency and freshness), MMKGs have some special quality issues that concern the images (e.g., wrong, missing or outdated facts), as shown in Table 6. Firstly, the image of some entity might be easily mixed with another when the two entities are closely related. Pluvianus aegyptius is a kind of bird that has a symbiosis with crocodiles, so we always get a picture of both the crocodile and the bird when searching for

it. Secondly, the images of a more famous entity may easily appear in the entity grounding results of its closely-related entities. The Wandering Earth is written by the famous Chinese science fiction writer *Liu Cixin*. While searching for this book, we always get a picture of his another more famous book, named The Dark Forest. Thirdly, some abstract concepts' visual features are not clear enough. For example, visual features of the arrogance are unfixed, so we always get some completely irrelevant pictures.

### 5.3 Efficiency

Efficiency is always a non-negligible issue when building a large-scale KG. The efficiency problem of constructing an MMKG is more striking, since the extra complexity of processing multimedia data needs to be considered. For example, it takes NEIL [22] around 350K CPU hours to collect 400K visual instances for 2273 objects, while in a typical KG we need to ground billions of instances. The scalability of the existing solutions in building MMKGs will be greatly challenged. If the grounding objective is video data, the scalability issue might be amplified.

Besides the construction of MMKG, the online application of MMKG also needs to carefully address the efficiency issue since the MMKG needs to serve applications in real-time. The solution's efficiency is crucial for online MMKG-based applications.

## 6 CONCLUSION

We are the first to thoroughly survey the existing work on MMKGs constructed by texts and images. We systematically review the existing work on MMKG construction and application. We compare mainstream MMKGs in terms of what they contain and how they construct. We analyze different solutions' strengths and weaknesses in MMKG construction and applications. We not only point out some potential opportunities with the existing tasks in both MMKG construction and application, but also list some promising future directions with the construction and application of MMKGs.

## REFERENCES

[1] C. Matuszek, M. Witbrock, J. Cabral, and J. DeOliveira, "An introduction to the syntax and content of cyc," *UMBC Computer Science and Electrical Engineering Department Collection*, 2006.

[2] H. Liu and P. Singh, "Conceptnet—a practical commonsense reasoning tool-kit," *BT technology journal*, 2004.

[3] G. A. Miller, "Wordnet: a lexical database for english," *Communications of the ACM*, 1995.

[4] R. Navigli and S. P. Ponzetto, "Babelnet: Building a very large multilingual semantic network," in *Proc. of ACL*, 2010.

[5] K. Bollacker, C. Evans, P. Paritosh, T. Sturge, and J. Taylor, "Freebase: a collaboratively created graph database for structuring human knowledge," in *Proc. of SIGMOD*, 2008.

[6] S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak, and Z. Ives, "Dbpedia: A nucleus for a web of open data," in *The semantic web*, 2007.

[7] F. M. Suchanek, G. Kasneci, and G. Weikum, "Yago: a core of semantic knowledge," in *Proc. of WWW*, 2007.

[8] D. Vrandečić and M. Krötzsch, "Wikidata: a free collaborative knowledgebase," *Communications of the ACM*, 2014.

[9] B. Xu, Y. Xu, J. Liang, C. Xie, B. Liang, W. Cui, and Y. Xiao, "Cn-dbpedia: A never-ending chinese knowledge extraction system," in *International Conference on Industrial, Engineering and Other Applications of Applied Intelligent Systems*, 2017.

[10] W. Wu, H. Li, H. Wang, and K. Q. Zhu, "Probase: A probabilistic taxonomy for text understanding," in *Proc. of SIGMOD*, 2012.

[11] M. Wick and B. Vatant, "The geonames geographical database," *Available from World Wide Web: http://geonames. org*, 2012.

[12] S. Harnad, "The symbol grounding problem," *Physica D: Nonlinear Phenomena*, 1990.

[13] ——, "Symbol grounding problem," *Encyclopedia of cognitive science*, 2003.

[14] L. Steels, "The symbol grounding problem has been solved. so what's next," *Symbols and embodiment: Debates on meaning and cognition*, 2008.

[15] E. M. Bender and A. Koller, "Climbing towards nlu: On meaning, form, and understanding in the age of data," in *Proc. of ACL*, 2020.

[16] B. Shi, L. Ji, P. Lu, Z. Niu, and N. Duan, "Knowledge aware semantic concept expansion for image-text matching," in *Proc. of IJCAI*, 2019.

[17] Y. Niu, K. Tang, H. Zhang, Z. Lu, X.-S. Hua, and J.-R. Wen, "Counterfactual vqa: A cause-effect look at language bias," in *Proc. of CVPR*, 2021.

[18] W. Zhao, Y. Hu, H. Wang, X. Wu, and J. Luo, "Boosting entity-aware image captioning with multi-modal knowledge graph," *arXiv preprint arXiv:2107.11970*, 2021.

[19] R. Sun, X. Cao, Y. Zhao, J. Wan, K. Zhou, F. Zhang, Z. Wang, and K. Zheng, "Multi-modal knowledge graphs for recommender systems," in *Proc. of CIKM*, 2020.

[20] G. Xu, H. Chen, F.-L. Li, F. Sun, Y. Shi, Z. Zeng, W. Zhou, Z. Zhao, and J. Zhang, "Alime mkg: A multi-modal knowledge graph for live-streaming e-commerce," in *Proc. of CIKM*, 2021.

[21] M. Li, A. Zareian, Y. Lin, X. Pan, S. Whitehead, B. Chen, B. Wu, H. Ji, S.-F. Chang, C. Voss *et al.*, "Gaia: A fine-grained multimedia knowledge extraction system," in *Proc. of ACL*, 2020.

[22] X. Chen, A. Shrivastava, and A. Gupta, "Neil: Extracting visual knowledge from web data," in *Proc. of ICCV*, 2013.

[23] H. Wen, Y. Lin, T. Lai, X. Pan, S. Li, X. Lin, B. Zhou, M. Li, H. Wang, H. Zhang *et al.*, "Resin: A dockerized schema-guided cross-document cross-lingual cross-media information extraction and event tracking system," in *Proc. of NAACL*, 2021.

[24] Y. Ma, Z. Wang, M. Li, Y. Cao, M. Chen, X. Li, W. Sun, K. Deng, K. Wang, A. Sun *et al.*, "Mmekg: Multi-modal event knowledge graph towards universal representation across modalities," in *Proc. of ACL*, 2022.

[25] S. Ferrada, B. Bustos, and A. Hogan, "Imgpedia: a linked dataset with content-based analysis of wikimedia images," in *Proc. of ISWC*, 2017.

[26] G. Vaidya, D. Kontokostas, M. Knuth, J. Lehmann, and S. Hellmann, "Dbpedia commons: structured multimedia metadata from the wikimedia commons," in *Proc. of ISWC*, 2015.

[27] D. Oñoro-Rubio, M. Niepert, A. García-Durán, R. González, and R. J. López-Sastre, "Answering visual-relational queries in web-extracted knowledge graphs," *Proc. of AKBC*, 2019.

[28] Y. Liu, H. Li, A. Garcia-Duran, M. Niepert, D. Onoro-Rubio, and D. S. Rosenblum, "Mmkg: Multi-modal knowledge graphs," in *European Semantic Web Conference*, 2019.

[29] M. Wang, H. Wang, G. Qi, and Q. Zheng, "Richpedia: a large-scale, comprehensive multi-modal knowledge graph," *Big Data Research*, 2020.

[30] H. Alberts, T. Huang, Y. Deshpande, Y. Liu, K. Cho, C. Vania, and I. Calixto, "Visualsem: a high-quality knowledge graph for vision and language," *arXiv preprint arXiv:2008.09150*, 2020.

[31] T. Baltrušaitis, C. Ahuja, and L.-P. Morency, "Multimodal machine learning: A survey and taxonomy," *IEEE transactions on pattern analysis and machine intelligence*, 2018.

[32] M. Cornia, L. Baraldi, and R. Cucchiara, "Show, control and tell: A framework for generating controllable and grounded captions," in *Proc. of CVPR*, 2019.

[33] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh, "Vqa: Visual question answering," in *Proc. of ICCV*, 2015.

[34] A. Karpathy, A. Joulin, and L. Fei-Fei, "Deep fragment embeddings for bidirectional image sentence mapping," *Proc. of*, 2014.

[35] S. K. D'mello and J. Kory, "A review and meta-analysis of multimodal affect detection systems," *ACM computing surveys (CSUR)*, 2015.

[36] R. Kiros, R. Salakhutdinov, and R. S. Zemel, "Unifying visual-semantic embeddings with multimodal neural language models," *arXiv preprint arXiv:1411.2539*, 2014.

[37] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio, "Show, attend and tell: Neural image caption generation with visual attention," in *Proc. of ICML*, 2015.

[38] D. Ramachandram and G. W. Taylor, "Deep multimodal learning: A survey on recent advances and trends," *IEEE signal processing magazine*, 2017.

[39] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang, "Bottom-up and top-down attention for image captioning and visual question answering," in *Proc. of CVPR*, 2018.

[40] T. Zhu, Y. Wang, H. Li, Y. Wu, X. He, and B. Zhou, "Multimodal joint attribute prediction and value extraction for e-commerce product," *arXiv preprint arXiv:2009.07162*, 2020.

[41] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," *arXiv preprint arXiv:2103.00020*, 2021.

[42] Y.-C. Chen, L. Li, L. Yu, A. E. Kholy, F. Ahmed, Z. Gan, Y. Cheng, and J. Liu, "Uniter: Learning universal image-text representations," *arXiv preprint arXiv:1909.11740*, 2019.

[43] L. H. Li, P. Zhang, H. Zhang, J. Yang, C. Li, Y. Zhong, L. Wang, L. Yuan, L. Zhang, J.-N. Hwang *et al.*, "Grounded language-image pre-training," in *Proc. of CVPR*, 2022.

[44] L. Yao, R. Huang, L. Hou, G. Lu, M. Niu, H. Xu, X. Liang, Z. Li, X. Jiang, and C. Xu, "Filip: Fine-grained interactive language-image pre-training," *arXiv preprint arXiv:2111.07783*, 2021.

[45] X. Li, X. Yin, C. Li, P. Zhang, X. Hu, L. Zhang, L. Wang, H. Hu, L. Dong, F. Wei *et al.*, "Oscar: Object-semantics aligned pre-training for vision-language tasks," in *Proc. of ECCV*, 2020.

[46] F. Yu, J. Tang, W. Yin, Y. Sun, H. Tian, H. Wu, and H. Wang, "Ernie-vil: Knowledge enhanced vision-language representations through scene graphs," in *Proc. of AAAI*, 2021.

[47] Y. Cui, Z. Yu, C. Wang, Z. Zhao, J. Zhang, M. Wang, and J. Yu, "Rosita: Enhancing vision-and-language semantic alignments via cross-and intra-modal knowledge integration," in *Proc. of ACM MM*, 2021.

[48] A. Mogadala, U. Bista, L. Xie, and A. Rettinger, "Describing natural images containing novel objects with knowledge guided assitance," *arXiv preprint arXiv:1710.06303*, 2017.

[49] K. Marino, M. Rastegari, A. Farhadi, and R. Mottaghi, "Ok-vqa: A visual question answering benchmark requiring external knowledge," in *Proc. of CVPR*, 2019.

[50] Q. Zhang, J. Fu, X. Liu, and X. Huang, "Adaptive co-attention network for named entity recognition in tweets," in *Proc. of AAAI*, 2018.

[51] K. Marino, X. Chen, D. Parikh, A. Gupta, and M. Rohrbach, "Krisp: Integrating implicit and symbolic knowledge for open-domain knowledge-based vqa," in *Proc. of CVPR*, 2021.

[52] P. Perona, "Vision of a visipedia," *Proceedings of the IEEE*, 2010.

[53] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma *et al.*, "Visual genome: Connecting language and vision using crowdsourced dense image annotations," *IJCV*, 2017.

[54] M. Warren and P. J. Hayes, "Bounding ambiguity: Experiences with an image annotation system." in *SAD/CrowdBias@ HCOMP*, 2018.

[55] M. Warren, D. A. Shamma, and P. J. Hayes, "Knowledge engineering with image data in real-world settings." in *Proc. of AAAI*, 2021.

[56] K. Chen, J. Gao, and R. Nevatia, "Knowledge aided consistency for weakly supervised phrase grounding," in *Proc. of CVPR*, 2018.

[57] V. Ramanishka, A. Das, J. Zhang, and K. Saenko, "Top-down visual saliency guided by captions," in *Proc. of CVPR*, 2017.

[58] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *Proc. of ECCV*, 2014.

[59] P. Young, A. Lai, M. Hodosh, and J. Hockenmaier, "From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions," *Transactions of the Association for Computational Linguistics*, 2014.

[60] B. A. Plummer, L. Wang, C. M. Cervantes, J. C. Caicedo, J. Hockenmaier, and S. Lazebnik, "Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models," in *Proc. of ICCV*, 2015.

[61] A. Kuznetsova, H. Rom, N. Alldrin, J. Uijlings, I. Krasin, J. Pont-Tuset, S. Kamali, S. Popov, M. Malloci, A. Kolesnikov, T. Duerig, and V. Ferrari, "The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale," *IJCV*, 2020.

[62] M. Wang, M. Azab, N. Kojima, R. Mihalcea, and J. Deng, "Structured matching for phrase localization," in *Proc. of ECCV*, 2016.

[63] W. Samek, G. Montavon, A. Vedaldi, L. K. Hansen, and K.-R. Müller, *Explainable AI: interpreting, explaining and visualizing deep learning*. Springer Nature, 2019.

[64] S. Hooker, D. Erhan, P.-J. Kindermans, and B. Kim, "A benchmark for interpretability methods in deep neural networks," *Proc. of NeurIPS*, 2019.

[65] P.-J. Kindermans, S. Hooker, J. Adebayo, M. Alber, K. T. Schütt, S. Dähne, D. Erhan, and B. Kim, "The (un) reliability of saliency methods," in *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, 2019.

[66] J. Bastings and K. Filippova, "The elephant in the interpretability room: Why use attention as explanation when we have saliency methods?" *arXiv preprint arXiv:2010.05607*, 2020.

[67] H. Akbari, S. Karaman, S. Bhargava, B. Chen, C. Vondrick, and S.-F. Chang, "Multi-level multimodal common semantic space for image-phrase grounding," in *Proc. of CVPR*, 2019.

[68] M. Li, A. Zareian, Q. Zeng, S. Whitehead, D. Lu, H. Ji, and S.-F. Chang, "Cross-media structured common space for multimedia event extraction," *arXiv preprint arXiv:2005.02472*, 2020.

[69] J. Zhang, S. A. Bargal, Z. Lin, J. Brandt, X. Shen, and S. Sclaroff, "Top-down neural attention by excitation backprop," *IJCV*, 2018.

[70] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, "Emerging properties in self-supervised vision transformers," *arXiv preprint arXiv:2104.14294*, 2021.

[71] P. Wang, A. Yang, R. Men, J. Lin, S. Bai, Z. Li, J. Ma, C. Zhou, J. Zhou, and H. Yang, "Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework," in *Proc. of ICML*, 2022.

[72] W. Kim, B. Son, and I. Kim, "Vilt: Vision-and-language transformer without convolution or region supervision," *arXiv preprint arXiv:2102.03334*, 2021.

[73] J. Hessel, A. Holtzman, M. Forbes, R. L. Bras, and Y. Choi, "Clipscore: A reference-free evaluation metric for image captioning," *arXiv preprint arXiv:2104.08718*, 2021.

[74] X. Zhang, X. Sun, C. Xie, and B. Lun, "From vision to content: Construction of domain-specific multi-modal knowledge graph," *IEEE Access*, 2019.

[75] G. Kulkarni, V. Premraj, V. Ordonez, S. Dhar, S. Li, Y. Choi, A. C. Berg, and T. L. Berg, "Babytalk: Understanding and generating simple image descriptions," *IEEE transactions on pattern analysis and machine intelligence*, 2013.

[76] D. Elliott and F. Keller, "Image description using visual dependency representations," in *Proc. of EMNLP*, 2013.

[77] B. Yao and L. Fei-Fei, "Grouplet: A structured image representation for recognizing human and object interactions," in *Proc. of CVPR*, 2010.

[78] ——, "Modeling mutual context of object and human pose in human-object interaction activities," in *Proc. of CVPR*, 2010.

[79] S. Maji, L. Bourdev, and J. Malik, "Action recognition from a distributed representation of pose and appearance," in *CVPR 2011*, 2011.

[80] S. Antol, C. L. Zitnick, and D. Parikh, "Zero-shot learning via visual abstraction," in *Proc. of ECCV*, 2014.

[81] R. Zellers, M. Yatskar, S. Thomson, and Y. Choi, "Neural motifs: Scene graph parsing with global context," in *Proc. of CVPR*, 2018.

[82] C. Lu, R. Krishna, M. Bernstein, and L. Fei-Fei, "Visual relationship detection with language priors," in *Proc. of ECCV*, 2016.

[83] B. Dai, Y. Zhang, and D. Lin, "Detecting visual relationships with deep relational networks," in *Proc. of CVPR*, 2017.

[84] H. Zhang, Z. Kyaw, S.-F. Chang, and T.-S. Chua, "Visual translation embedding network for visual relation detection," in *Proc. of CVPR*, 2017.

[85] D. Xu, Y. Zhu, C. B. Choy, and L. Fei-Fei, "Scene graph generation by iterative message passing," in *Proc. of CVPR*, 2017.

[86] J. Yang, J. Lu, S. Lee, D. Batra, and D. Parikh, "Graph r-cnn for scene graph generation," in *Proc. of ECCV*, 2018.

[87] R. Yu, A. Li, V. I. Morariu, and L. S. Davis, "Visual relationship detection with internal and external linguistic knowledge distillation," in *Proc. of ICCV*, 2017.

[88] K. Tang, H. Zhang, B. Wu, W. Luo, and W. Liu, "Learning to compose dynamic tree structures for visual contexts," in *Proc. of CVPR*, 2019.

[89] J. Zhang, K. J. Shih, A. Elgammal, A. Tao, and B. Catanzaro, "Graphical contrastive losses for scene graph parsing," in *Proc. of CVPR*, 2019.

[90] W. Wang, R. Liu, M. Wang, S. Wang, X. Chang, and Y. Chen, "Memory-based network for scene graph with unbalanced relations," in *Proc. of ACM MM*, 2020.

[91] W. Wang, M. Wang, S. Wang, G. Long, L. Yao, G. Qi, and Y. Chen, "One-shot learning for long-tail visual relation detection," in *Proc. of AAAI*, 2020.

[92] J. Zhang, Y. Kalantidis, M. Rohrbach, M. Paluri, A. Elgammal, and M. Elhoseiny, "Large-scale visual relationship understanding," in *Proc. of AAAI*, 2019.

[93] K. Tang, Y. Niu, J. Huang, J. Shi, and H. Zhang, "Unbiased scene graph generation from biased training," in *Proc. of CVPR*, 2020.

[94] M. Baumgartner, L. Rossetto, and A. Bernstein, "Towards using semantic-web technologies for multi-modal knowledge graph construction," in *Proc. of ACM MM*, 2020.

[95] Y.-L. Li, L. Xu, X. Liu, X. Huang, Y. Xu, M. Chen, Z. Ma, S. Wang, H.-S. Fang, and C. Lu, "Hake: Human activity knowledge engine," *arXiv preprint arXiv:1904.06539*, 2019.

[96] T. Zhang, S. Whitehead, H. Zhang, H. Li, J. Ellis, L. Huang, W. Liu, H. Ji, and S.-F. Chang, "Improving event extraction via multimodal integration," in *Proc. of ACM MM*, 2017.

[97] B. Chen, X. Lin, C. Thomas, M. Li, S. Yoshida, L. Chum, H. Ji, and S.-F. Chang, "Joint multimedia event extraction from video and article," *arXiv preprint arXiv:2109.12776*, 2021.

[98] Y. Xiong, K. Zhu, D. Lin, and X. Tang, "Recognize complex events from static images by fusing deep channels," in *Proc. of CVPR*, 2015.

[99] M. Yatskar, L. Zettlemoyer, and A. Farhadi, "Situation recognition: Visual semantic role labeling for image understanding," in *Proc. of CVPR*, 2016.

[100] S. Pratt, M. Yatskar, L. Weihs, A. Farhadi, and A. Kembhavi, "Grounded situation recognition," in *Proc. of ECCV*, 2020.

[101] H. Li, J. G. Ellis, H. Ji, and S.-F. Chang, "Event specific multimodal pattern mining for knowledge base construction," in *Proc. of ACM MM*, 2016.

[102] L. Banarescu, C. Bonial, S. Cai, M. Georgescu, K. Griffitt, U. Hermjakob, K. Knight, P. Koehn, M. Palmer, and N. Schneider, "Abstract meaning representation for sembanking," in *Proceedings of the 7th linguistic annotation workshop and interoperability with discourse*, 2013.

[103] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proc. of CVPR*, 2009.

[104] K. Yang, K. Qinami, L. Fei-Fei, J. Deng, and O. Russakovsky, "Towards fairer datasets: Filtering and balancing the distribution of the people subtree in the imagenet hierarchy," in *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, 2020.

[105] A. Torralba, R. Fergus, and W. T. Freeman, "80 million tiny images: A large data set for nonparametric object and scene recognition," *IEEE transactions on pattern analysis and machine intelligence*, 2008.

[106] Y. Yan, G. Liu, S. Wang, J. Zhang, and K. Zheng, "Graph-based clustering and ranking for diversified image search," *Multimedia Systems*, 2017.

[107] X. Jiang, A. Li, J. Liang, B. Liu, R. Xie, W. Wu, Z. Li, and Y. Xiao, "Visualizable or non-visualizable? exploring the visualizability of concepts in multi-modal knowledge graph," in *Proc. of DASFAA*, 2022.

[108] Q. Mei-Bin, Z. Jun-Jun, J. Ping, and J. Jian-Guo, "Representative image selection from image dataset," *ACTA AUTOMATICA SINICA*, 2014.

[109] H. Yu, Z.-H. Deng, Y. Yang, and T. Xiong, "A joint optimization model for image summarization based on image content and tags," in *Proc. of AAAI*, 2014.

[110] Y. Wang, L. Zhu, and X. Qian, "Social image retrieval based on topic diversity," *Multimedia Tools and Applications*, 2021.

[111] R. H. Van Leuken, L. Garcia, X. Olivares, and R. van Zwol, "Visual diversification of image search results," in *Proc. of WWW*, 2009.

[112] M. Wang, K. Yang, X.-S. Hua, and H.-J. Zhang, "Towards a relevant and diverse search of social images," *IEEE Transactions on Multimedia*, 2010.

[113] J. Xueyao, L. Weichen, L. Jingping, L. Zhixu, and X. Yanghua, "Entity image collection based on multi-modality pattern transfer," *Computer Engineering*, 2022.

[114] T. Deselaers, T. Gass, P. Dreuw, and H. Ney, "Jointly optimising relevance and diversity in image retrieval," in *Proceedings of the ACM international conference on image and video retrieval*, 2009.

[115] Z. Ji, Y. Su, Y. Pang, and X. Qu, "Diversifying the image relevance reranking with absorbing random walks," in *2011 Sixth International Conference on Image and Graphics*, 2011.

[116] R. Raguram and S. Lazebnik, "Computing iconic summaries of general visual concepts," in *Proc. of CVPR*, 2008.

[117] M. Elhoseiny, S. Cohen, W. Chang, B. Price, and A. Elgammal, "Sherlock: Scalable fact learning in images," in *Proc. of AAAI*, 2017.

[118] Y. Guo, J. Chen, H. Zhang, and Y.-G. Jiang, "Visual relations augmented cross-modal retrieval," in *Proc. of ICMR*, 2020.

[119] S. Wang, R. Wang, Z. Yao, S. Shan, and X. Chen, "Cross-modal scene graph matching for relationship-aware image-text retrieval," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2020.

[120] C. Zheng, Z. Wu, J. Feng, Z. Fu, and Y. Cai, "Mnre: A challenge multimodal dataset for neural relation extraction with visual evidence in social media posts," in *ICME*, 2021.

[121] C. Zheng, J. Feng, Z. Fu, Y. Cai, Q. Li, and T. Wang, "Multimodal relation extraction with efficient graph alignment," in *Proc. of ACM MM*, 2021.

[122] S. Moon, L. Neves, and V. Carvalho, "Zeroshot multimodal named entity disambiguation for noisy social media posts," in *Proc. of ACL*, 2018.

[123] L. Zhang, Z. Li, and Q. Yang, "Attention-based multimodal entity linking with high-quality images," in *Proc. of DASFAA*, 2021.

[124] X. Wang, J. Ye, Z. Li, J. Tian, Y. Jiang, R. Wang, M. Yan, J. Zhang, and Y. Xiao, "Cat-mner: Multimodal named entity recognition with knowledge-refined cross-modal attention," in *ICME*, 2022.

[125] D. A. Hudson and C. D. Manning, "Gqa: A new dataset for real-world visual reasoning and compositional question answering," in *Proc. of CVPR*, 2019.

[126] P. Wang, Q. Wu, C. Shen, A. Dick, and A. van den Hengel, "Fvqa: Fact-based visual question answering," *IEEE transactions on pattern analysis and machine intelligence*, 2018.

[127] S. Shah, A. Mishra, N. Yadati, and P. P. Talukdar, "Kvqa: Knowledge-aware visual question answering," in *Proc. of AAAI*, 2019.

[128] P. Wang, Q. Wu, C. Shen, A. v. d. Hengel, and A. Dick, "Explicit knowledge-based reasoning for visual question answering," *arXiv preprint arXiv:1511.02570*, 2015.

[129] T.-S. Chua, J. Tang, R. Hong, H. Li, Z. Luo, and Y. Zheng, "Nus-wide: a real-world web image database from national university of singapore," in *Proceedings of the ACM international conference on image and video retrieval*, 2009.

[130] S. Guadarrama, N. Krishnamoorthy, G. Malkarnenkar, S. Venugopalan, R. Mooney, T. Darrell, and K. Saenko, "Youtube2text: Recognizing and describing arbitrary activities using semantic hierarchies and zero-shot recognition," in *Proc. of ICCV*, 2013.

[131] A. Tran, A. Mathews, and L. Xie, "Transform and tell: Entity-aware news image captioning," in *Proc. of CVPR*, 2020.

[132] T.-H. Huang, F. Ferraro, N. Mostafazadeh, I. Misra, A. Agrawal, J. Devlin, R. Girshick, X. He, P. Kohli, D. Batra *et al.*, "Visual storytelling," in *Proc. of NAACL*, 2016.

[133] F. M. Harper and J. A. Konstan, "The movielens datasets: History and context," *Acm transactions on interactive intelligent systems (tiis)*, 2015.

[134] A. Uyar and F. M. Aliyu, "Evaluating search features of google knowledge graph and bing satori: entity types, list searches and query interfaces," *Online Information Review*, 2015.

[135] Y. Huang, M. Li, and Y. Wu, "Kkbox's music recommendation."

[136] J. Xiong, G. Liu, Y. Liu, and M. Liu, "Oracle bone inscriptions information processing based on multi-modal knowledge graph," *Computers & Electrical Engineering*, 2021.

[137] N. Li, Q. Shen, R. Song, Y. Chi, and H. Xu, "Medukg: A deep-learning-based approach for multi-modal educational knowledge graph construction," *Information*, 2022.

[138] A. V. Kannan, D. Fradkin, I. Akrotirianakis, T. Kulahcioglu, A. Canedo, A. Roy, S.-Y. Yu, M. Arnav, and M. A. Al Faruque, "Multimodal knowledge graph for deep learning papers and code," in *Proc. of CIKM*, 2020.

[139] C. Deng, Y. Jia, H. Xu, C. Zhang, J. Tang, L. Fu, W. Zhang, H. Zhang, X. Wang, and C. Zhou, "Gakg: A multimodal geoscience academic knowledge graph," in *Proc. of CIKM*, 2021.

[140] P. Bloem, X. Wilcke, L. v. Berkel, and V. d. Boer, "kgbench: A collection of knowledge graph datasets for evaluating relational and multimodal machine learning," in *European Semantic Web Conference*, 2021.

[141] W. Zheng, L. Yan, C. Gou, Z.-C. Zhang, J. J. Zhang, M. Hu, and F.-Y. Wang, "Pay attention to doctor-patient dialogues: Multi-modal knowledge graph attention image-text embedding for covid-19 diagnosis," *Information Fusion*, 2021.

[142] A. Carlson, J. Betteridge, B. Kisiel, B. Settles, E. R. Hruschka, and T. M. Mitchell, "Toward an architecture for never-ending language learning," in *Proc. of AAAI*, 2010.

[143] M. Nickel, V. Tresp, and H. Kriegel, "A three-way model for collective learning on multi-relational data," in *Proc. of ICML*, 2011.

[144] Z. Zhang, Z. Li, H. Liu, and N. N. Xiong, "Multi-scale dynamic convolutional network for knowledge graph embedding," *IEEE Trans. Knowl. Data Eng.*, 2022.

[145] P. Pezeshkpour, L. Chen, and S. Singh, "Embedding multimodal relational data for knowledge base completion," 2018.

[146] A. Bordes, N. Usunier, A. Garcia-Duran, J. Weston, and O. Yakhnenko, "Translating embeddings for modeling multi-relational data," *Proc. of NeurIPS*, 2013.

[147] Z. Wang, J. Zhang, J. Feng, and Z. Chen, "Knowledge graph embedding by translating on hyperplanes," in *Proc. of AAAI*, 2014.

[148] Y. Lin, Z. Liu, M. Sun, Y. Liu, and X. Zhu, "Learning entity and relation embeddings for knowledge graph completion," in *Proc. of AAAI*, 2015.

[149] G. Ji, S. He, L. Xu, K. Liu, and J. Zhao, "Knowledge graph embedding via dynamic mapping matrix," in *Proc. of ACL*, 2015.

[150] H. Moussely-Sergieh, T. Botschen, I. Gurevych, and S. Roth, "A multimodal translation-based approach for knowledge graph representation learning," in *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, 2018.

[151] A. Rettinger, "Towards holistic concept representations: Embedding relational knowledge, visual attributes, and distributional word semantics," 2017.

[152] X. Chen, N. Zhang, L. Li, S. Deng, C. Tan, C. Xu, F. Huang, L. Si, and H. Chen, "Hybrid transformer with multi-level fusion for multimodal knowledge graph completion," *arXiv preprint arXiv:2205.02357*, 2022.

[153] R. Xie, Z. Liu, H. Luan, and M. Sun, "Image-embodied knowledge representation learning," 2017.

[154] A. Rossi, D. Barbosa, D. Firmani, A. Matinata, and P. Merialdo, "Knowledge graph embedding for link prediction: A comparative analysis," *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 2021.

[155] W. Wilcke, P. Bloem, V. de Boer, R. van t Veer, and F. van Harmelen, "End-to-end entity classification on multimodal knowledge graphs," *arXiv preprint arXiv:2003.12383*, 2020.

[156] H. Guo, J. Tang, W. Zeng, X. Zhao, and L. Liu, "Multi-modal entity alignment in hyperbolic space," *Neurocomputing*, 2021.

[157] L. Chen, Z. Li, Y. Wang, T. Xu, Z. Wang, and E. Chen, "Mmea: Entity alignment for multi-modal knowledge graph," in *Proc. of KSEM*, 2020.

[158] A. Garcia-Duran and M. Niepert, "Kblrn: End-to-end learning of knowledge base representations with latent, relational, and numerical features," *Proc. of UAI*, 2018.

[159] D. Chen, Z. Li, B. Gu, and Z. Chen, "Multimodal named entity recognition with image attributes and image knowledge," in *Proc. of DASFAA*, 2021.

[160] O. Adjali, R. Besançon, O. Ferret, H. Le Borgne, and B. Grau, "Multimodal entity linking for tweets," *Advances in Information Retrieval*, 2020.

[161] S. Moon, L. Neves, and V. Carvalho, "Multimodal named entity disambiguation for noisy social media posts," in *Proc. of ACL*, 2018.

[162] J. Yu, Z. Zhu, Y. Wang, W. Zhang, Y. Hu, and J. Tan, "Cross-modal knowledge reasoning for knowledge-based visual question answering," *Pattern Recognition*, 2020.

[163] S. Bhakthavatsalam, K. Richardson, N. Tandon, and P. Clark, "Do dogs have whiskers? a new knowledge base of haspart relations," *arXiv preprint arXiv:2006.07510*, 2020.

[164] L. Ma, Z. Lu, L. Shang, and H. Li, "Multimodal convolutional neural networks for matching image and sentence," in *Proc. of ICCV*, 2015.

[165] Y. Huang, Q. Wu, C. Song, and L. Wang, "Learning semantic concepts and order for image and sentence matching," in *Proc. of CVPR*, 2018.

[166] K.-H. Lee, X. Chen, G. Hua, H. Hu, and X. He, "Stacked cross attention for image-text matching," in *Proc. of ECCV*, 2018.

[167] L. Ma, W. Jiang, Z. Jie, Y.-G. Jiang, and W. Liu, "Matching image and sentence with multi-faceted representations," *IEEE Transactions on Circuits and Systems for Video Technology*, 2019.

[168] Z. Zheng, L. Zheng, M. Garrett, Y. Yang, M. Xu, and Y.-D. Shen, "Dual-path convolutional image-text embeddings with instance loss," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 2020.

[169] Y. Huang, Q. Wu, W. Wang, and L. Wang, "Image and sentence matching via semantic concepts and order learning," *IEEE transactions on pattern analysis and machine intelligence*, 2018.

[170] W. Wang, Y. Huang, and L. Wang, "Language-driven temporal activity localization: A semantic matching reinforcement learning model," in *Proc. of CVPR*, 2019.

[171] H. Wang, Y. Zhang, Z. Ji, Y. Pang, and L. Ma, "Consensus-aware visual-semantic embedding for image-text matching," in *Proc. of ECCV*, 2020.

[172] C. Chaudhary, P. Goyal, D. N. Prasad, and Y.-P. P. Chen, "Enhancing the quality of image tagging using a visio-textual knowledge base," *IEEE Transactions on Multimedia*, 2019.

[173] J. Hou, X. Wu, Y. Qi, W. Zhao, J. Luo, and Y. Jia, "Relational reasoning using prior knowledge for visual captioning," *arXiv preprint arXiv:1906.01290*, 2019.

[174] C.-C. Hsu, Z.-Y. Chen, C.-Y. Hsu, C.-C. Li, T.-Y. Lin, T.-H. Huang, and L.-W. Ku, "Knowledge-enriched visual storytelling," in *Proc. of AAAI*, 2020.

[175] J. Bobadilla, F. Ortega, A. Hernando, and A. Gutiérrez, "Recommender systems survey," *Knowledge-based systems*, 2013.

[176] X. Shen, B. Yi, H. Liu, W. Zhang, Z. Zhang, S. Liu, and N. Xiong, "Deep variational matrix factorization with knowledge embedding for recommendation system," *IEEE Trans. Knowl. Data Eng.*, 2021.

[177] X. Su and T. M. Khoshgoftaar, "A survey of collaborative filtering techniques," *Advances in artificial intelligence*, 2009.

[178] F. Zhang, N. J. Yuan, D. Lian, X. Xie, and W.-Y. Ma, "Collaborative knowledge base embedding for recommender systems," in *Proc. of KDD*, 2016.

[179] S. Tao, R. Qiu, Y. Ping, and H. Ma, "Multi-modal knowledge-aware reinforcement learning network for explainable recommendation," *Knowledge-Based Systems*, 2021.

**Xiangru Zhu** is a Ph.D. student with the School of Computer Science at Fudan University, China. Her research interests include multi-modal knowledge graph and vision-language pre-trained model.



**Zhixu Li** is a professor with the School of Computer Science at Fudan University, China. He used to be a professor at Soochow University between 2014 and 2021. He received his Ph.D. degree in Computer Science from the University of Queensland in 2013. His main research interests are Data & Knowledge Engineering, and Cognitive Intelligence, and he is particularly interested in Multi-modal Knowledge Graph and Cross-Modal Cognitive Intelligence.



**Xiaodan Wang** is a Master student with the School of Computer Science at Fudan University, China. Her research interests include image-text retrieval and multi-modal knowledge graphs construction.



**Xueyao Jiang** received her Master degree in computer science from Fudan University in 2022. Her research interests include multi-modal knowledge graphs construction and application.



**Penglei Sun** is a master student with the School of Computer Science at Fudan University, China. His research interests include scenario-driven multi-modal knowledge graph construction and application.



**Xuwu Wang** is a Ph.D. student with the School of Computer Science at Fudan University, China. Her research interests mainly focus on entity recognition, entity linking and multi-modal knowledge acquisition and application. She has already published several papers on ACL, ICME, DASFAA etc.



**Yanghua Xiao** is a professor of computer science at Fudan University. He is the director of Knowledge Works Lab, Fudan University. He got his PHD degree in software theory from Fudan University, Shanghai, China, in 2009. He is one of young 973 scientists. His research interest includes big data management and mining, graph database, knowledge graph. He was a visiting professor of Human Genome Sequencing Center at Baylor College Medicine, and visiting researcher of Microsoft Research Asia.



**Nicholas Jing Yuan** (Senior Member, IEEE) is currently a General Manager of AI Services, the Chief Scientist, and the Director of the Language and Speech Innovation Lab, Huawei Cloud. He has published more than 60 papers in top-tier conferences and journals, including several best paper awards such as SIGKDD (2016, 2018), ICDM (2013), and SIGSPATIAL (2010). His research work has been featured by influential media such as MIT Technology Review many times, and was reported directly to Bill Gates (Founder of Microsoft) by himself. He served regularly as program committee members in top tier conferences such as SIGKDD, WWW, ACL, and AAAI. He is a Senior Member of ACM and CCF.