

Are a Thousand Words Better Than a Single Picture?

Beyond Images - A Framework for Multi-Modal Knowledge Graph Dataset Enrichment

Supplementary Material

1 Image Embedding Methods in MMKGs

Semantically Ambiguous Images. Current MMKG models typically embed images as vectors and combine them with embeddings from other modalities (e.g., text) to create richer entity representations. These embedding methods generally follow two approaches:

- Global feature extraction: Methods like Convolutional Neural Networks (CNNs) generate fixed-size global feature vectors for entire images (Li et al., 2022). While efficient for large-scale datasets, they often fail to capture fine-grained details.
- Local feature extraction: Approaches such as Vision Transformers (ViTs) divide images into patches and embed each patch individually, enabling finer-grained feature extraction and improved alignment with text (Dosovitskiy et al., 2020). However, these methods are computationally intensive and heavily reliant on image quality and alignment effectiveness.

Challenges in Handling Specific Image Types.

However, both approaches face challenges when processing certain types of image, where standard embedding methods fail to capture essential semantic features.

- Sparse-Semantic Images (e.g., symbolic logos): These images contain limited visual information, often featuring simple geometric shapes or elements. While they may carry critical domain knowledge, existing models struggle to extract distinctive embeddings, reducing their effectiveness (Su et al., 2024; Wang et al., 2020).
- Rich-Semantic Images (e.g., abstract paintings): These images are visually and semantically complex, including intricate scenes, interactions, or artistic expressions. Current

embedding methods often struggle to fully capture these semantic relationships, leading to significant information loss (Wilber et al., 2017).

2 Standardizing and Aligning Original Images

To generate textual descriptions based on images, we require the original images from the datasets. However, many existing works only provide image embeddings (vectors) without the raw images. For those that do provide raw images, naming conventions for image files vary significantly. Some use Wikidata URLs, others use DBpedia URLs, and some rely on YAGO entity names. This inconsistency, especially with YAGO entity names that often include special characters such as \, /, or :, creates challenges in aligning image with entity names. Many operating systems are unable to handle filenames containing such characters, further complicating the alignment process and subsequent experiments.

To address this, we first standardize the naming conventions for raw images in the datasets. Concretely, we align all entities using their Wikidata IDs (QIDs). The QIDs consist only of alphanumeric characters, which are compatible with all operating systems, facilitating future reproduction and extensions. Additionally, QIDs serve as a bridge between entities and their Wikipedia pages, enabling us to download supplementary images and metadata (e.g., timestamps) from Wikipedia for dataset enrichment. Details are summarized in Table 1.

MKG-W. We found that original images for different entities were stored in folders named after the entities, but many special characters (e.g., \, /, or :) were missing. Additionally, the image filenames within these folders lacked any recognizable pattern. To address this, we used the dataset's

Table 1: Overview of three public MMKG datasets, summarizing key statistics including the number of entities, relations, dataset splits, and image attributes. The table details original images, newly downloaded images, average images per entity, and images with timestamps.

	MKG-W	MKG-Y	DB15K
Entity	15,000	15,000	12,842
Relation	169	28	279
Train	34,196	21,310	79,222
Validation	4,276	2,665	9,902
Test	4,274	2,663	9,904
Text	14,123	12,305	9,078
Original Images			
Total Img	27,841	42,242	603,435
Avg Img	3.00	3.00	53.35
Img w/ Timestamp	0	0	0
Entity w/ Img	9,285	14,099	11,311
New Images			
Total Img	81,323	56,646	176,858
Avg Img	5.81	4.23	14.58
Img w/ Timestamp	55,317	39,281	124,721
Entity w/ Img	14,002	14,388	12,130

provided mapping file, which links DBpedia URLs to Wikidata URLs, to identify the corresponding entity names and QIDs for each entity.

Next, we removed all special characters from both the extracted entity names and the folder names containing the original images to facilitate matching. Once the matching was complete, we had the following information for each sample: entity name, QID, and original images. Finally, we renamed all images using the format `qid_idx` and consolidated them into a single folder for use in subsequent experiments.

MKG-Y. We followed a similar process as in MKG-W (see above). The original images were stored in folders named after the entities, but the filenames lacked a consistent naming convention. Unlike MKG-W, the original dataset did not provide a mapping file between DBpedia and Wikidata URLs. However, it did include a mapping between DBpedia URLs and sample names.

Using the DBpedia URLs, we accessed the corresponding DBpedia pages and leveraged `sameAs` links to locate the corresponding Wikidata pages and obtain the QIDs. We then matched the folders containing raw images to their respective entities and renamed the images using the `qid_idx` format. Finally, all renamed images were consolidated into a single folder for subsequent use.

DB15K. The original paper (Liu et al., 2019) did not provide downloadable images, only image embeddings and URLs for the images. As a result,

we re-downloaded the images using the provided links. Each sample had 100 links from Google Images, approximately 35 from Bing Images, and 50 from Yahoo Image Search. According to the original paper, the top 20 images from each search engine (for a total of 60 images per entity) should be downloaded.

However, some links were no longer valid. To ensure fairness in reproducing the results, we sequentially downloaded up to 20 images from each search engine. If fewer than 20 valid images were available, we continued downloading from subsequent links until 20 images were obtained per search engine, maintaining the original dataset’s image count of 60 per sample.

After downloading the images, we used the DBpedia URLs to access the DBpedia pages, followed `sameAs` links to locate the corresponding Wikidata pages, and obtained the QIDs for each sample. Finally, we renamed the images using the `qid_idx` format and consolidated them into a single folder for subsequent experiments.

3 Our Datasets Structure

During processing, images in each batch were encoded as tensors using a processor and then passed through the model to generate textual descriptions. The generated descriptions and their corresponding filenames were saved to an output file, providing semantically meaningful textual data for subsequent MMKG tasks. A summary of the final output files is provided in Table 2.

Table 2: Each image in the dataset includes unique identifiers, source URLs, metadata (e.g., date, author), and BLIP-2-generated textual descriptions.

Key	Description
<code>id</code>	Unique identifier for each image
<code>page_url</code>	URL of the Wikipedia page
<code>image_url</code>	URL of the image file
<code>table_data</code>	Metadata of the image
- <code>Description</code>	Brief description of the image
- <code>Date</code>	Date associated with the image
- <code>Author</code>	Author or creator of the image
- <code>Formatted_Date</code>	Standardized date format
<code>image_blip2_detail</code>	Detailed textual description

The dataset provides detailed information for each image, such as URLs, metadata, and automatically generated textual descriptions. Each entry consists of a unique identifier (`id`), links to the corresponding Wikipedia page (`page_url`) and image file (`image_url`), metadata extracted from the

image (table_data), and textual descriptions generated by BLIP-2 (image_blip2_detail). The metadata includes various attributes such as date (Formatted_Date), author, and resolution, which are crucial for MMKG research.

4 Implementation Details

Our experiments use the default hyperparameters for each Baseline model to ensure fair comparisons. All experiments were conducted on a Linux server equipped with a single NVIDIA H100 GPU. To generate textual descriptions from images, we used the BLIP-2 model, with English as the output language. The maximum generated text length is limited to 100 words. On average, each image generated 20 words, ranging from 15 to 25 words. The generated text is then embedded into vectors using BERT-base-uncased.

5 Computational Cost

Our method works directly with existing images in MMKG datasets and achieves 3-4% improvement (Figure 3 in the paper), validating our main hypothesis(RQ1). Retrieving new images from the internet is optional and was only used for the ablation study(RQ2). While external images can bring additional gains (Table 3 and Figure 3 in the paper), they are not required for the method to be effective.

Image retrieval and captioning is a one-time pre-processing step, similar to standard MMKG setups. As detailed in the supplementary (Sec. 2&4), generating descriptions for 100K images takes 1 hour (MKG-W/Y) or 8 hours (DB15K) on a single H100 GPU. The output text files are <200MB and reusable.

Once generated, this enriched data can be reused across models and tasks without extra cost. As shown in Figure 3 in the paper, training time increases by only 7-30 minutes, while performance improves, demonstrating a favorable cost-benefit trade-off.

6 Image-to-Text Generation Models

“*Blip2-flan-t5-xxl*” efficiently bridges images and text by integrating a frozen image encoder with a frozen language model, connected through a lightweight Querying Transformer. This design is particularly well-suited for our task, as it effectively translates visual features into meaningful textual descriptions while minimizing computational

overhead and eliminating the need for extensive retraining.

“*Git-large-coco*” model adopts a generative Transformer architecture that directly converts CLIP image tokens into textual descriptions. Fine-tuned on the COCO image captioning dataset, it is well-suited for generating coherent and relevant captions. Its end-to-end training allows it to extract visual details and express them in natural language effectively.

“*Llava-v1.5-7b*” model integrates a pre-trained CLIP vision encoder with a large language model connected through a lightweight projection layer. Fine-tuned on visual instruction data, it generates detailed, context-aware descriptions conditioned on both the image and the user prompt. It leverages the generative and reasoning capabilities of the language model to produce high-quality image-grounded text.

7 Main Result

The main results are shown in Table 3, which summarizes the link prediction performance of four models (MMRNS, MyGO, NativeE, and AdaMF) across three datasets (MKG-W, MKG-Y, and DB15K) under different settings.

Model performance is evaluated using rank-based metrics, including Mean Reciprocal Rank (MRR) and Hits@ K ($K = 1, 3, 10$). MRR calculates the average of the reciprocal ranks of the correct answers in the predicted ranking list, while Hits@ K measures the proportion of correct answers appearing within the top K predictions. Both metrics are commonly used in evaluating link prediction tasks, with higher scores indicating better model performance.

The first row for each model presents the experimental results on the original datasets, as reproduced from the original papers. “*G*” stands for *Generate*, referring to our framework that generates textual descriptions from images. “*o*” indicates that the textual descriptions were generated from the original images provided in the dataset, while “*n*” means that the descriptions were generated from images automatically downloaded using our *Beyond Images* framework. *Improvement* represents the percentage increase ($\text{Boost} = \frac{\text{Our Result} - \text{Baseline Result}}{\text{Baseline Result}}$) in performance of the enriched datasets compared to the original datasets.

Table 3: Link prediction results of four models across three datasets. “*D*” represents entity descriptions, “*I*” denotes image embeddings, “*G(o)*” refers to textual descriptions generated from original images, and “*G(n)*” corresponds to textual descriptions from newly downloaded images. “*H@n*” stands for “Hits at *n*.” The “*Improvement* (↑%)” indicates the performance gain of the best-performing model (highlighted in bold) over the Baseline model.

Model	MKG-W				MKG-Y				DB15K			
	MRR	H@1	H@3	H@10	MRR	H@1	H@3	H@10	MRR	H@1	H@3	H@10
MMRNS	35.03	28.59	37.49	47.47	35.93	30.53	39.07	45.47	32.68	23.01	37.86	51.01
MMRNS (D+I+G(o))	35.73	29.65	38.37	48.69	36.59	31.78	40.19	46.43	33.57	24.04	39.13	52.71
Improvement (↑%)	1.99%	3.70%	2.35%	2.57%	1.84%	4.11%	2.87%	2.11%	2.74%	4.49%	3.35%	3.34%
MMRNS (D+I+G(n))	36.13	29.93	38.58	49.02	36.93	31.96	40.33	46.58	33.37	23.78	39.02	52.40
Improvement (↑%)	3.13%	4.68%	2.91%	3.26%	2.79%	4.69%	3.21%	2.44%	2.11%	3.36%	3.06%	2.72%
MMRNS (D+I+G(o+n))	36.26	30.08	38.70	49.19	37.03	32.12	40.46	46.70	33.67	24.16	39.27	52.89
Improvement (↑%)	3.50%	5.20%	3.21%	3.62%	3.07%	5.21%	3.56%	2.71%	3.04%	5.01%	3.71%	3.69%
MyGO	36.10	29.78	38.54	47.75	38.51	33.39	39.03	47.87	37.72	30.08	41.26	52.21
MyGO (D+I+G(o))	37.19	30.85	39.65	48.75	39.63	34.73	39.88	48.90	38.84	31.53	42.37	53.74
Improvement (↑%)	3.01%	3.61%	2.88%	2.10%	2.91%	4.01%	2.17%	2.14%	2.97%	4.81%	2.69%	2.92%
MyGO (D+I+G(n))	37.28	31.26	39.74	49.18	39.83	35.07	40.20	49.22	38.77	31.23	42.30	53.24
Improvement (↑%)	3.28%	4.97%	3.12%	2.99%	3.42%	5.03%	3.01%	2.81%	2.78%	3.81%	2.53%	1.97%
MyGO (D+I+G(o+n))	37.42	31.42	39.88	49.35	39.97	35.26	40.32	49.37	38.97	31.69	42.49	53.92
Improvement (↑%)	3.66%	5.51%	3.48%	3.35%	3.80%	5.60%	3.31%	3.13%	3.31%	5.36%	2.99%	3.27%
NativE	36.58	29.56	39.65	48.94	39.04	34.79	40.89	46.18	37.16	28.01	41.36	54.13
NativE (D+I+G(o))	37.37	30.56	40.44	49.93	39.63	35.95	41.93	47.03	38.68	28.83	42.48	55.11
Improvement (↑%)	2.16%	3.38%	1.98%	2.02%	1.52%	3.33%	2.55%	1.83%	4.10%	2.92%	2.72%	1.81%
NativE (D+I+G(n))	37.57	30.68	40.85	50.27	39.75	36.12	42.11	47.43	38.30	28.77	42.35	55.03
Improvement (↑%)	2.72%	3.80%	3.02%	2.72%	1.81%	3.83%	2.99%	2.71%	3.08%	2.72%	2.41%	1.66%
NativE (D+I+G(o+n))	37.69	30.80	40.97	50.41	39.83	36.27	42.25	47.56	38.84	28.92	42.61	55.22
Improvement (↑%)	3.02%	4.21%	3.32%	3.01%	2.01%	4.25%	3.32%	2.98%	4.52%	3.25%	3.02%	2.02%
AdaMF	35.85	29.04	39.01	48.42	38.57	34.34	40.59	45.76	35.14	25.30	41.11	52.92
AdaMF (D+I+G(o))	36.92	30.16	39.78	49.34	39.79	35.37	41.45	46.41	36.20	26.24	42.29	54.35
Improvement (↑%)	2.98%	3.84%	1.96%	1.90%	3.16%	2.99%	2.12%	1.43%	3.02%	3.71%	2.87%	2.71%
AdaMF (D+I+G(n))	37.20	30.35	39.77	49.73	40.05	35.86	41.89	46.78	35.85	26.08	42.13	54.24
Improvement (↑%)	3.77%	4.51%	1.94%	2.70%	3.84%	4.44%	3.20%	2.23%	2.01%	3.08%	2.49%	2.49%
AdaMF (D+I+G(o+n))	37.36	30.50	39.85	49.88	40.21	36.04	42.02	46.88	36.32	26.34	42.43	54.51
Improvement (↑%)	4.21%	5.02%	2.15%	3.02%	4.25%	4.95%	3.52%	2.46%	3.35%	4.12%	3.20%	3.00%

8 Case Analysis: Boosting Performance with Textual Descriptions

8.1 Example 1

All images are shown in Figure 1. Triple: (*Hot Sauce Committee Part Two*, *performer*, *Beastie Boys*). Images (a) and (b) correspond to the head entity *Hot Sauce Committee Part Two*, while images (c) - (h) represent the tail entity *Beastie Boys*.

Triple: (Hot Sauce Committee Part Two, performer, Beastie Boys)

QID: (Q1933719, P175, Q214039)

Head entity’s rank: correct head entity’s rank improved from 13,680 to 1,330.

Tail entity’s rank: correct tail entity’s rank improved from 11,435 to 4,628.

8.2 Example 2

All images are shown in Figure 2. Triple: (*Her Harem*, *cast member*, *Carroll Baker*). Images (a) - (c) correspond to the head entity *Her Harem*, while images (d) - (m) represent the tail entity *Carroll*

Baker.

Triple: (Her Harem, cast member, Carroll Baker)

QID: (Q3819142, P161, Q233891)

Head entity’s rank: correct head entity’s rank improved from 10,177 to 8,611.

Tail entity’s rank: correct tail entity’s rank improved from 571 to 72.

8.3 Example 3

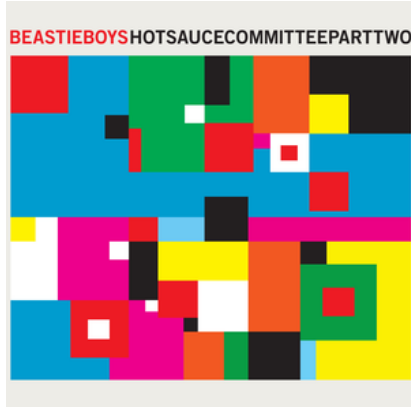
All images are shown in Figure 3. Triple: (*World (The Price of Love)*, *performer*, *New Order*). Images (a) correspond to the head entity *World (The Price of Love)*, while images (b) - (f) represent the tail entity *New Order*.

Triple: (World (The Price of Love), performer, New Order)

QID: (Q8035321, P175, Q214990)

Head entity’s rank: correct head entity’s rank improved from 12,528 to 2,622.

Tail entity’s rank: correct tail entity’s rank improved from 10,185 to 2,591.



(a) Q1933719_1: “The cover of beastboys hot sauce committee part two”.



(b) Q1933719_2: “The scene shows a group of men walking on a bridge”.



(c) Q214039_1: “three men are leaning on a stair railing”.



(d) Q214039_2: “The logo for beastie boys is shown in black and white”.



(e) Q214039_3: “two men are standing on stage with a microphone”.



(f) Q214039_4: “two men in black jackets are on stage singing”.



(g) Q214039_5: “a man in a red suit and hat is singing on stage”.



(h) Q214039_6: “a man in a suit and tie singing”.

Figure 1: Triple: (*Hot Sauce Committee Part Two*, *performer*, *Beastie Boys*). Images (a) and (b) correspond to the head entity *Hot Sauce Committee Part Two*, while images (c) - (h) represent the tail entity *Beastie Boys*.



(a) Q3819142_1: “the poster for the movie *her harem*”.
 (b) Q3819142_2: “the italian flag is shown on a clapperboard”.
 (c) Q3819142_3: “two theatrical masks on a clapperboard”.
 (d) Q233891_1: “a black and white photo of a woman with long blonde hair”.
 (e) Q233891_2: “a black background with a white tv screen”.

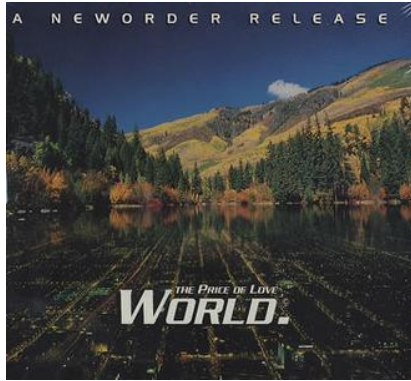


(f) Q233891_3: “a woman in a black and white photo”.
 (g) Q233891_4: “a man and woman in western attire sit on a horse”.
 (h) Q233891_5: “a woman in a striped top sits on a bench”.
 (i) Q233891_6: “a woman in a fur coat sits on a white fur rug”.



(j) Q233891_7: “a woman is standing in a shower”.
 (k) Q233891_8: “the scene shows a man and woman talking to each other”.
 (l) Q233891_9: “a woman in a white dress is standing on a stage in front of a large ship”.
 (m) Q233891_10: “a star on the hollywood walk of fame for carroll baker”.

Figure 2: Triple: (*Her Harem*, cast member, *Carroll Baker*). Images (a) - (c) correspond to the head entity *Her Harem*, while images (d) - (m) represent the tail entity *Carroll Baker*.



(a) Q8035321_1: “the cover of the world album”.



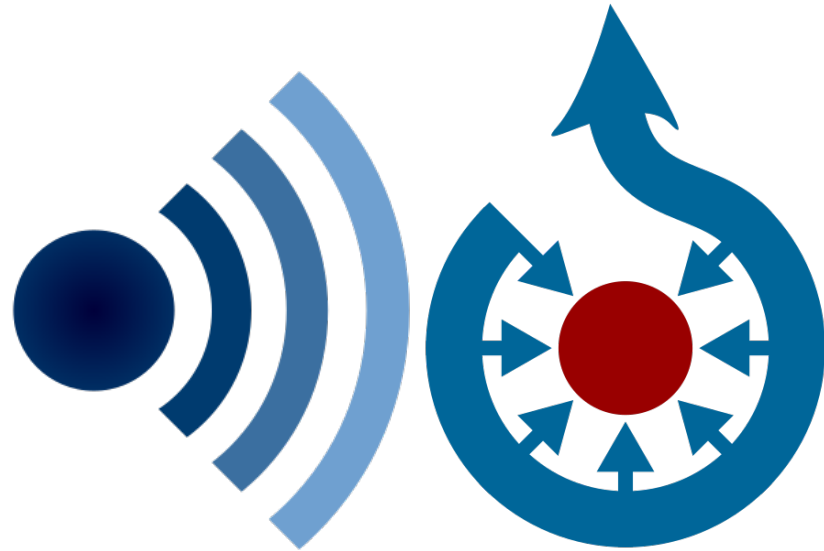
(b) Q214990_1: “four black and white photos of four men”.



(c) Q214990_2: “a group of men are on stage with guitars and drums”.



(d) Q214990_3: “a band is performing on stage with a large screen behind them”.



(e) Q214990_4: “a blue and white wave symbol”.

(f) Q214990_5: “a blue and red logo with arrows pointing in different directions”.

Figure 3: Triple: (*World (The Price of Love)*, *performer*, *New Order*). Images (a) correspond to the head entity *World (The Price of Love)*, while images (b) - (f) represent the tail entity *New Order*.

References

- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2020. [An image is worth 16x16 words: Transformers for image recognition at scale](#). *ArXiv*, abs/2010.11929.
- Zewen Li, Fan Liu, Wenjie Yang, Shouheng Peng, and Jun Zhou. 2022. [A survey of convolutional neural networks: Analysis, applications, and prospects](#). *IEEE Transactions on Neural Networks and Learning Systems*, 33(12):6999–7019.
- Ye Liu, Hui Li, Alberto Garcia-Duran, Mathias Niepert, Daniel Onoro-Rubio, and David S. Rosenblum. 2019. [Mmkg: Multi-modal knowledge graphs](#). In *The Semantic Web: 16th International Conference, ESWC 2019, Portorož, Slovenia, June 2–6, 2019, Proceedings*, page 459–474, Berlin, Heidelberg. Springer-Verlag.
- Taoyu Su, Xinghua Zhang, Jiawei Sheng, Zhenyu Zhang, and Tingwen Liu. 2024. [Loginmea: Local-to-global interaction network for multi-modal entity alignment](#). In *ECAI 2024*, pages 1173–1180. IOS Press.
- Jing Wang, Weiqing Min, Sujuan Hou, Shengnan Ma, Yuanjie Zheng, Haishuai Wang, and Shuqiang Jiang. 2020. [Logo-2k+: A large-scale logo dataset for scalable logo classification](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(04):6194–6201.
- Michael J Wilber, Chen Fang, Hailin Jin, Aaron Hertzmann, John Collomosse, and Serge Belongie. 2017. [Bam! the behance artistic media dataset for recognition beyond photography](#). In *Proceedings of the IEEE international conference on computer vision*, pages 1202–1211.