# Building a Large-scale Multimodal Knowledge Base System
# for Answering Visual Queries

Yuke Zhu     Ce Zhang     Christopher Ré     Li Fei-Fei
Computer Science Department, Stanford University
{yukez,czhang,chrismre,feifeili}@cs.stanford.edu

## Abstract

*The complexity of the visual world creates significant challenges for comprehensive visual understanding. In spite of recent successes in visual recognition, today's vision systems would still struggle to deal with visual queries that require a deeper reasoning. We propose a knowledge base (KB) framework to handle an assortment of visual queries, without the need to train new classifiers for new tasks. Building such a large-scale multimodal KB presents a major challenge of scalability. We cast a large-scale MRF into a KB representation, incorporating visual, textual and structured data, as well as their diverse relations. We introduce a scalable knowledge base construction system that is capable of building a KB with half billion variables and millions of parameters in a few hours. Our system achieves competitive results compared to purpose-built models on standard recognition and retrieval tasks, while exhibiting greater flexibility in answering richer visual queries.*

## 1. Introduction

Type the following query in Google (i.e., a search engine) – "names of universities in Manhattan". The returned list of answers is often sensible. But try this one – "names of universities with computer science PhD program in Manhattan". The answers are far from satisfying. Both questions are perfectly clear to most humans, but current NLP-based algorithms still fail to perform well for more complex queries. In vision, we see a similar pattern. Much progress has been made in tasks such as classification and detection on single objects (e.g., Fig. 1(a)). But real-world vision applications might require more diverse and heterogeneous querying needs (e.g., Fig. 1(b)). The traditional classification-based methods would struggle in such tasks.
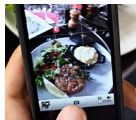
Towards the goal of scaling up the large-scale, diverse and heterogeneous visual querying tasks, a handful of recent papers [7, 59] have suggested to cast the visual recognition tasks into a framework that enables more heteroge-



(a) **Find me pictures of a dog.**

(b)

**Q: Where can I find similar cuisines in downtown Chicago?**

**Answers:**

U. U. Grill
Chicago, IL 60642

S. C. Steak House
Chicago, IL 60657

**Q: Find photos of me sea kayaking last Halloween in my photo album.**

**Answers**: Saturday, October 31

Figure 1: Although a classification-based method might be sufficient to find images of a dog in query (a). It would struggle for queries in real-world applications. To answer the queries in (b), we need to fuse visual information with metadata for joint reasoning. We propose a visual knowledge base framework to perform different types of visual tasks without training new classifiers. Our framework allows one to express this complex task with a single query.

neous reasoning and inference. A major benefit of doing so is to avoid training a new set of classifiers every time a new type of questions arises. We approach this problem by building a large-scale multimodal *knowledge base* (KB), where we answer visual queries (like the ones in Fig. 1(b)) by evaluating probabilistic KB queries.

A KB can often be viewed as a large-scale graph structure that connects different entities with their relations [38, 58]. In NLP, some early promising results have been shown by encoding entity and relation information in text-based KBs, e.g., Freebase [3] and IBM Waston's Jeopardy system [13]. In vision, there is now a small but growing amount of attention in building visual KBs. In NEIL [7], Chen et al. have shown the benefit of using contextual relations between scenes, objects and attributes to improve scene classi-

fication and object detection. However, its testing scenario is limited on recognition-based tasks; while it lacks a coherent inference model to extend to richer high-level tasks without training new classifiers. Zhu et al. [59] have shown how to build a Markov Logic KB for affordance reasoning. However, their testing scenario is limited by its small data size and the discrete representation. Our paper is particularly inspired by these two works [7, 59], but focuses on addressing the following two key challenges.

First, **answering a variety of heterogeneous visual queries without re-training**. In real-world vision applications, the space of possible queries is huge (even infinite). It is impossible to retrain classifiers for every type of queries. Our system demonstrates its ability to perform reasoning and inference on an assortment of visual querying tasks, ranging from scene classification, image search to real-world application queries, without the need to train new classifiers for new tasks. We formalize answering these visual queries as computing the marginal probabilities of the joint probability model (Sec. 5). The key technique is to express visual queries in a logical form that can be answered from the visual KB in a principled inference method. We qualitatively evaluate our KB model in answering application queries like the ones in Fig. 1 (Sec. 6.1). We then perform quantitative evaluations on the recognition tasks (Sec. 6.2) and retrieval tasks (Sec. 6.3) respectively using the SUN dataset [50]. Our system achieves competitive results compared to the classification-based baseline models, while exhibiting greater flexibility in answering a variety of visual queries.

Second, **learning with large-scale multimodal data**. To build such a scalable KB, the model needs to perform joint learning and inference on a large amount of images, text and structured data, especially by using both discrete and continuous variables. Existing text-based KB representations [15, 38, 58] fail to incorporate continuous visual features in a probabilistic framework, which hinders us from expressing richer multimodal data. In vision, MRFs have been widely used as a probabilistic framework to model joint distributions among multimodal variables. We cast a MRF model into a KB representation to accommodate a mixture of discrete and continuous variables in a joint probability model. While MRFs have been widely used in a variety of vision tasks [9, 26, 27, 46], applying them to a large-scale KB framework means that we need to conquer the challenge of scalable learning and inference. We build a scalable visual KB construction system by leveraging database techniques, high-speed sampling [55] and first-order methods [35]. We are able to build a KB with half billion variables and four million parameters, which is four orders of magnitude larger than Zhu et al. [59] while using half of its training time.

## 2. Previous Work

**Joint Models in Vision**    A series of context models have leveraged MRFs in various vision tasks, such as image segmentation [16, 27, 33], object recognition [9, 26], object detection [46], pose and activity recognition [52] and other recognition tasks [20, 36]. Similarly, the family of And-Or graph models [47, 56] focus on parsing images and videos into a hierarchical structure. In this work, we use an MRF representation for joint learning and inference of our data, casting MRF models into modern KB systems. In particular, we address the scalability challenge of large-scale MRF learning with our knowledge base construction system.

**Learning with Vision and Language**    Previous work on joint learning with vision and language abounds [23, 30, 41, 42, 60]. Image and video captioning has recently become a popular task, where the goal is to generate a short text description for images and videos [8, 11, 21, 29, 44, 48, 51]. It is followed by visual question answering [1, 14, 31, 32, 53], which aims at answering natural language questions based on image content. Both captioning and question answering tasks perform on a single image and produce NLP outputs. Our system offers one single, coherent framework that can perform joint learning and inference on one or multiple images as well as metadata in textual and other forms.

**Knowledge Bases**    Most KB work in the database and NLP communities focuses on organizing and retrieving only textual information in a structured representation [3, 13, 28, 58]. Although a few large-scale KBs [3, 12] have made attempts to incorporate visual information, they simply cache the visual contents and link them to text via hyperlinks. In vision, a series of work has focused on extracting relational knowledge from visual data [5, 39, 60]. Chen et al. [7], Divvala et al. [10] and Zhu et al. [59] have recently proposed KB-based frameworks for visual recognition tasks. However, they all lack an inference framework to deal with more diverse types of visual queries. PhotoRecall [25] proposed a pre-defined knowledge structure to retrieve photos from text queries. In contrast, our system allows for new KB structures and offers the flexibility of answering richer types of queries.

## 3. A Joint Probability Model: Casting a Large-Scale MRF into a KB System

Our first task is to build a system that can efficiently learn a KB given a large amount of multimodal information, such as images, metadata, textual labels, and structured labels. Towards a real-world, large-scale system like this, the challenges are two-fold. First, our learning system must allow for a coherent probabilistic representation of both discrete and continuous variables to accommodate the heterogeneity of the data. Second, we need to develop an efficient but principled learning and inference method that is capable of

large-scale computation. We address the first property in this section, and the second in Sec. 4.

### 3.1. The Knowledge Base System

A KB can be intuitively thought of as a graph of nodes connected by edges as in Fig. 2, where the nodes are called "entities" and the edges are called "relations". In vision, MRFs have been widely used to represent such graph structures [20, 33, 36, 46]. Thus, we cast an MRF model as the KB representation, where entities are represented by variables and relations by edges between variables. This model provides an umbrella framework for answering visual queries, where we formalize query answering as evaluating marginals from the joint distribution (Sec. 5). In comparison to MLNs used in previous work [38, 59], this representation is more generic, allowing us to accommodate continuous random variables and real-valued factors. In practice, we use factor graphs [24, 49], a bipartite graph equivalence of an MRF. Factor graphs provide a simple graphical interpretation of the MRF model, resulting in ease of implementation for large-scale inference.

A factor graph has two types of nodes: *variables* and *factors*. A possible world is a particular assignment to every variable, denoted by $I$. We define the probability of a possible world $I$ to be proportional to a log-linear combination of factors. We assign different weights to factors, expressing their relative influence on the probability. Formally, we define the *partition function* $Z$ of a possible world $I$ as

$$Z[I] = \exp\left(\sum_{i=1}^{m} w_i f_i(I)\right) \qquad (1)$$

where $w_i$ is the weight of the $i$-th factor, $f_i(I)$ is the value of the $i$-th factor in possible world $I$, and $m$ is the total number of factors. The probability of a possible world is

$$\Pr[I; \mathbf{w}] = Z[I]\left(\sum_{I' \in \mathcal{I}} Z[I']\right)^{-1} \qquad (2)$$

where $\mathcal{I}$ is the set of all possible worlds, and $\mathbf{w}$ corresponds to the factor weights. In Fig. 2, each node corresponds to a variable; and each edge between nodes corresponds to a factor. We define all the factors used in our KB in Sec. 3.2.

Having defined the structure of the factor graph KB, our learning objective is to find the optimal weight

$$\mathbf{w}^* = \arg\min_{\mathbf{w}} -\sum_{I \in \mathcal{I}_E} \log \Pr[I; \mathbf{w}] + \lambda||\mathbf{w}||_2^2 \qquad (3)$$

where $\mathcal{I}_E$ is the set of possible worlds obtained from the training images and $\lambda$ is the regularization parameter. To optimize Eq. (3), we need to compute the stochastic gradient $\frac{\partial \Pr[I|\mathbf{w}]}{\partial \mathbf{w}}$. It is usually intractable to compute the analyt-
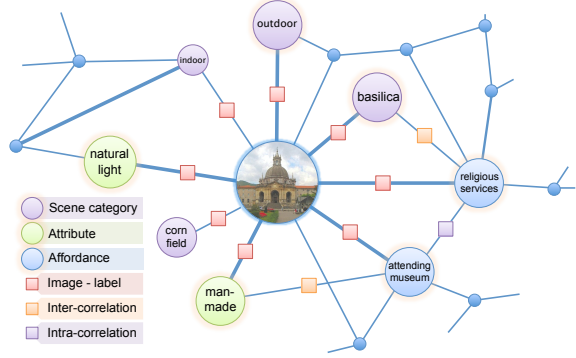


Figure 2: **A graphical illustration of a visual knowledge base (KB).** A visual KB contains both visual entities (e.g., scene images) and textual entities (e.g., semantic labels) interconnected by various types of edges characterizing their relations. The nodes and edges correspond to the variables and factors respectively in the factor graph. The colors indicate different node (edge) types.

ical gradients, as it involves the computation of an expectation over all possible words. We use the contrastive divergence scheme [19] to estimate the log-likelihood gradients. The gradient of the weight $w_i$ of the $i$-th factor (omitting regularization) is approximated by:

$$\nabla w_i \approx f_i(I') - f_i(I'') \qquad (4)$$

where $I'$ is a possible world sampled from the training data, and $I''$ is a possible world sampled under the distribution formed by the model (parameterized by $\mathbf{w}$). Gibbs sampling [6, 17] is used as the transition operator of the Markov chain. Intuitively, the first term in Eq. (4) increases the probability of training data; and the second term decreases the probability of samples generated by the model. In-depth studies on the estimated gradients of Eq. (4) can be found in the context of RBM training [4, 45]. We show in Sec. 4 that our system automatically creates a factor graph and learns the weights in a principled and scalable manner.

### 3.2. Data Sources for the Knowledge Base

We now describe the entities and relations in our KB, and the data sources that we will use to populate the KB. For our purposes, SUN [50] is a particularly useful dataset because of a) its diverse set of images, and b) the availability of a large number of category and attribute labels.

**Entities** can be thought of as descriptors of the images. In the factor graph depicted in Fig. 2, they are the nodes (variables) of the graph.

*Images* – are represented by their 4096-dimensional activations from the last fully-connected layer in a convolutional network [54]. In total, there are 59,709 images from the SUN dataset [50], where half are used for building the KB, and half for evaluation.

*Scene category labels* – indicate scene classes. In our experiments, we use 15 basic-level categories (e.g., workplace and transportation), and 298 fine-grained level categories (e.g., grotto and swamp) from SUN [50].

*Attribute labels* – characterize visual properties (e.g., material, layouts, lighting, etc.) of a scene. We use the SUN Attribute Dataset [37], which provides 102 attribute labels (e.g., glossy and warm).

*Affordance labels* – describe the functional properties of a scene, i.e., the actions that one can perform in a scene. We use a lexicon of 227 affordances (actions).[1] We conducted a large-scale online experiment to annotate the possibilities of the 227 actions for each scene category. We provide the list of affordances in Sec. C in the supplementary material.

**Relations** link entities (variables) to each other, as depicted by the squares on the edges in Fig. 2. The weights learned for the edges (factors) indicate the strength of the relations. We introduce three types of relations in our model.

*Image - label* – maps image features to semantic labels.

*Intra-correlations* – capture the co-occurrence between attribute-attribute and affordance-affordance pairs.

*Inter-correlations* – characterize correlations between two different types of labels (category - affordance, affordance - attribute, category - attribute and relations between categories from different levels).

The entities and relations in the KB are mapped to variables and factors in the factor graph. We represent the image entities as continuous variables, and the label entities as discrete variables. Each image is associated with hundreds of attribute and affordance labels. Together, this amounts to a KB of millions of entities. Table 1 summarizes some of the basic statistics of the KB that will be learned. This is two orders of magnitude larger than previous work [59] regarding the number of entities and relations. The large size of our dataset presents a significant challenge of scalability. In theory, an MRF can be arbitrarily large. However, its scalability is subject to the inefficiency of learning and inference. In addition, it is prohibitive to handcraft such a large-scale model from scratch. We, therefore, need a principled and scalable system for constructing the visual KB.

Table 1: **KB Dataset Statistics**

|  | Attributes | Affordances |
|---|---|---|
| **Lexicon size** | 102 | 227 |
| **# Total labels** | $1.34 \times 10^6$ | $1.36 \times 10^7$ |
| **# Positive labels** | $9.6 \times 10^5$ | $1.23 \times 10^6$ |
| **# Positive / image** | 6.7 | 13.7 |

---

[1]from the American Time Use Survey (ATUS) [40] sponsored by the Bureau of Labor Statistics, which catalogs the actions in daily lives and represents United States census data

# 4. Learning the Large-scale KB System

Given our goal towards learning a real-world, large-scale MRF-based KB system, the biggest challenge we need to address here is efficient learning and inference. A number of recent advances have been made in the database community to shed light on how to build a large-scale KB [12, 13, 34]. Our framework follows closely that of Niu et al. [34]. In addition to that, we address the challenge of learning with multimodal data. Our KB system and the data will be made available to the public.

## 4.1. Scalable Construction

There are three key steps to make the knowledge base construction (KBC) scalable: data pre-processing, factor graph generation and high-performance learning. Fig. 3 offers an overview of the KBC process illustrating these three steps, which are indicated by the boxes.

**Data Pre-processing** The first step (the first box in Fig. 3) is to pre-process raw data into a structured representation, in particular, as tables in a relational database. Each database table stores the entities of the same type (e.g., the Affordance table in Fig. 3(a)). It provides us access to database techniques such as SQL queries and parallel computing, important to achieve high scalability. We provide the database schema in Sec. A in the supplementary material.

**Factor Graph Generation** We represent the MRF model by a factor graph for the ease of implementation for scalable learning. The factor graph is generated from the database tables (the second box in Fig. 3). Each row in the database tables corresponds to a variable in the factor graph. For each training image, we construct a factor graph, where the variables (blue circles in Fig. 3(b)) are linked to their values in the database (dashed lines between Fig. 3(a) and (b)). We then define the factors on these variables. It is prohibitive to handcraft a large KB structure. Instead, we develop a declarative language that allows us to define the factors with a handful of human-readable rules. This language is a simple but powerful extension to previous work like MLNs [38] and PRMs [15], which enables us to specify relations between multimodal entities in logical conjunctions. We show an example rule in Fig. 3(b). This rule describes co-occurrence between affordance label `travel` and attribute label `sunny` on image `I1`. It evaluates to 1 if both labels are true and 0 otherwise. The KBC system parses this rule and creates a factor $f_k$ on these two variables. A weight $w_k$ is assigned to this factor and will be learned in the next step. The system creates a small factor graph for each of the training images. There is no edge between these graphs; however, the same factors in the graphs share the same weight (illustrated by the red squares in Fig. 3(c)). The weight sharing scheme is also specified in the declarative language. We provide a detailed explanation of the
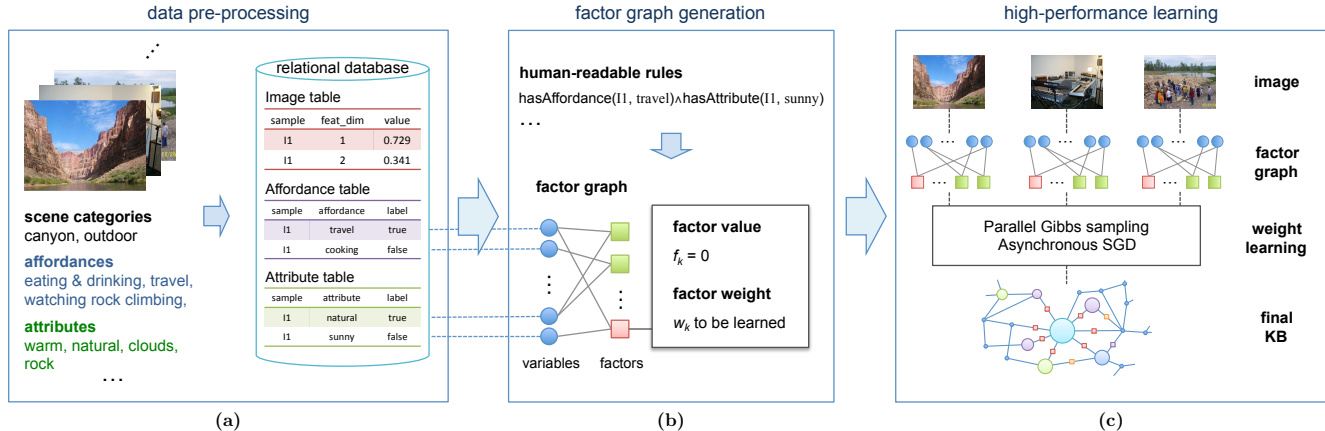
Figure 3: **An overview of the knowledge base construction pipeline.** We first process the images and text, converting them into a structured representation. We write human-readable rules to define the KB structure. The system automatically creates a factor graph by parsing the rules. We then adopt a scalable Gibbs sampler to learn the weights in the factor graph.

declarative language and a complete list of rules in Sec. A in the supplementary material.

**High-Performance Learning** Having defined the factor graph structures, our goal is to learn the factor weights efficiently. We use the learning method in Sec. 3.1 to find the optimal factor weights. We built a Gibbs sampler for high-performance learning and inference that is able to handle multimodal variables. Our system performs scalable Gibbs sampling based on careful system design and speedup techniques. On the system side, we implemented the Hogwild! model [35, 55] which can run asynchronous stochastic gradient descent while still guaranteeing convergence. The system runs in parallel, allowing the sampler to achieve a high efficiency. On average, our Gibbs sampler processes $8.2 \times 10^7$ variables per second. Finally this step produces a learned visual KB.

## 4.2. Learning Efficiency

The three steps (described in Sec. 4.1) together contribute to the high scalability of our KBC system. Table 2 shows that with this framework, we can build a KB four orders of magnitude larger regarding the number of variables and three orders of magnitude larger regarding model parameters compared to [59] (using Alchemy MLNs [38]), in half of the time. Fig. 4 further demonstrates that the learning time grows steadily as the KB size increases. The end-to-end construction finishes in 5.2 hours on the whole dataset (Sec. 3.2), indicating the potential to build larger-scale KBs in the future.

Table 2: **Statistics of the Visual KB Systems**

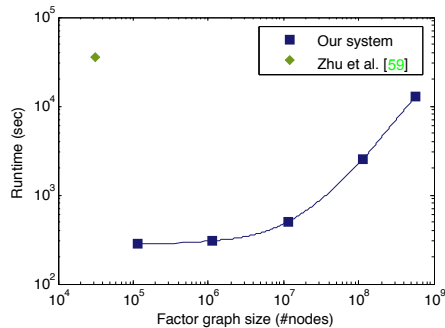|  | variables | parameters | runtime |
|---|---|---|---|
| **Zhu et al.** [59] | $3.15 \times 10^4$ | $5.06 \times 10^3$ | 10 hr |
| **Ours** | $5.76 \times 10^8$ | $4.19 \times 10^6$ | 5.2 hr |



Figure 4: **Efficiency of the knowledge base construction system.** The curve is plotted in log-log scale, where the $x$-axis is the number of nodes in the factor graph, and the $y$-axis is the runtime to construct the KB.

## 5. Visual Query Setup

As we have mentioned in the introduction, one advantage of using a KB system is its ability to handle rich and diverse types of visual queries without training new classifiers. Moreover, this inference is done in one joint model without step-wise filtering, treating images and other meta-data on an equal footing in learning and inference. From a user's perspective, the input to this system is a natural language question along with a set of one or more images. Similarly, the output is a mixture of images and text.

In practice, the space of possible queries is huge. It would be prohibitive to map each natural language question to the corresponding inference task in an ad-hoc manner. One solution is to reformulate the questions in a formal language [2], such as a probabilistic query language based on conjunctive queries [43]. This language allows us to express KB queries and to compute a ranked list of answers based on their marginal probabilities. We briefly describe how this works by an example query that retrieves

images of a sunny beach. This query is formed by a conjunction of two predicates (Boolean-valued functions) of `sceneCategory` and `hasAttribute`:

$$\texttt{sceneCategory}(i,\texttt{beach}) \land \texttt{hasAttribute}(i,\texttt{sunny})$$

Given such a query, our task is to find all possible images $i$ where both predicates are true – i.e., image $i$ comes from the scene category `beach` and has the attribute `sunny`. Following this example, more complex queries can be formed by joining several predicates together.[2]

Let $Q$ be a conjunctive query such as the one above. We compute a ranked list of answers (e.g., images of sunny beaches) based on their marginal probabilities. Formally, the marginal probability of a tuple $t$ (a list of variable assignments) being an answer to $Q$ is defined as:

$$\Pr[t \in Q] = \sum_{I \in \mathcal{I}} \mathbb{1}_{t \in Q(I)} \cdot \Pr[I; \mathbf{w}] \qquad (5)$$

where $\mathcal{I}$ and $\Pr[I; \mathbf{w}]$ are defined in Eq. (2), $\mathbb{1}$ is the indicator function, and $Q(I)$ is the set of variable assignments in the possible world $I$ under which $Q$ evaluates to true. We use the same Gibbs sampler as in Sec. 4.1 to estimate tuple marginals by sampling a collection of possible worlds and averaging the query values over these possible worlds. Each query evaluation produces a set of tuple-probability pairs $\{(t_1, p_1), (t_2, p_2), \ldots\}$, where we retrieve the top answers by sorting the pairs based on their probabilities in a descending order.

# 6. Experiments

Now that we have learned a large KB from multimodal data sources, and have established a probabilistic language to express visual queries, we can demonstrate how a KB can be useful in a number of querying tasks. To demonstrate the utility of our KB, we perform several types of evaluations that involve vision tasks with multimodal answers including images, text and metadata.

## 6.1. Answering Queries of Diverse Types

We start with a qualitative demonstration of using the KB to answer a wide variety of queries by performing joint inference on image appearance, as well as metadata like geolocations, timestamps, and business information.[3] Fig. 5 provides a few examples that depict the rich queries the system can handle. A user can ask the KB a question in natural language, such as "find me a modern looking mall

---

[2]In this work, we manually annotate the conjunctive queries from natural language questions. The mapping from sentences to logical forms is a well-studied problem in NLP [2] and orthogonal to our system.

[3]These metadata are either acquired from existing databases or automatically scraped online. Detailed descriptions of the experimental setups and the conjunctive queries (Sec. 5) for Fig. 5 are provided in Sec. B in the supplementary material.
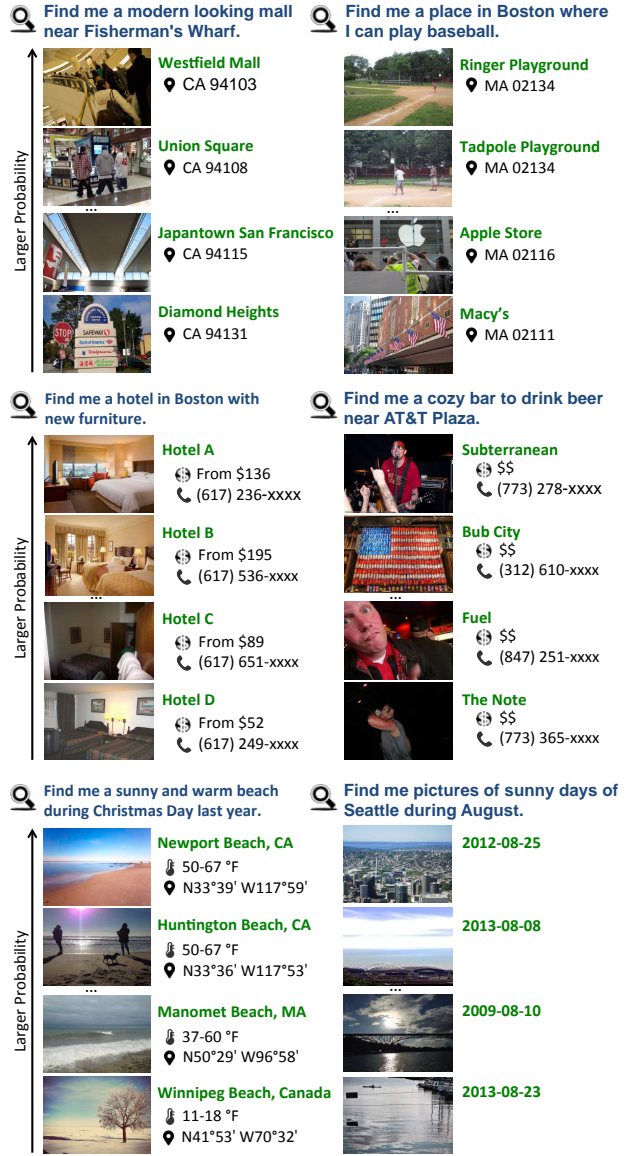


Figure 5: **Proof-of-concept queries in a query answering application.** We incorporate external data to enrich our knowledge base, and demonstrate its flexibility in answering real-world queries.

near Fisherman's Wharf." While the photos of the malls are not part of the training data in Sec. 3.2, our system is capable of linking the photo contents to other metadata, and is able to offer the names and locations of the shopping malls. Similarly in the second example "find me a place in Boston where I can play baseball", our system predicts the affordances from the appearances of the photos, and combines them with geolocation information to retrieve a list of places for playing baseball. In Fig. 5, the answers are shown in a ranked list by their marginal probabilities. Without a principled inference model, previous work such as NEILL[7] and LEVAN [10] cannot produce such probabilisitic outputs.

Table 3: **Performance of Scene Classification (in mAcc)**

| Method | Basic level | Fine-grained |
|---|---|---|
| CNN Fine-tuned [54] | 89.1 | 67.5 |
| Attribute-based model | 88.0 | 57.9 |
| Attributes + Features | 90.2 | 69.6 |
| KB - Affordances | 90.0 | 69.3 |
| KB - Attributes | 90.7 | 69.6 |
| KB - Full | **91.2** | **69.8** |

Table 4: **Performance of Scene Affordance Prediction**

| Method | mF1 | mAP |
|---|---|---|
| CNN Fine-tuned [54] | 81.6 | 74.2 |
| KB - Full | **82.6** | **75.7** |

## 6.2. Single-Image Query Answering

While our KB is designed for answering a wide range of queries, we can still evaluate how our system performs quantitatively in several standard visual recognition tasks without re-training. Based on the KB we have learned from data sources such as SUN (see Sec 3.2), we show two experiments for scene classification and affordance prediction. Both of these two tasks can be thought of as answering queries for a single image, where these queries can be expressed by a single predicate with the querying labels taken as random variables – i.e., sceneCategory(img, $c$) and hasAffordance(img, $a$). Our system outperforms the state-of-the-art baseline methods for each of these tasks.

For both experiments, we use the data in Sec. 3.2 for training and an evaluation set of 29,781 images from the same 298 categories of SUN [50] for testing. We measure *scene classification* by mean accuracy (mAcc) over classes [57]. SUN [50] provides two ways of classification: basic-level (15 categories) and fine-grained (298 categories). Table 3 provides a summary of the results, comparing our full model (KB - Full) with a number of different settings and state-of-the-art models. We describe the models used in Table 3 as follow:

- **CNN Fine-tuned** We fine-tuned a CNN [54] on a subset of SUN397 dataset [50] of 107,754 images. We train $\ell_2$-logistic regression classifiers on the activations from the last fully-connected layer. We also use this as image features for all the other baselines.

- **Attribute-based model** We predict the scene attributes and affordances from the CNN features, and use a binary vector of the predicted values as an intermediate feature. This is the strategy adopted by Zhu et al. [59] to discretize visual data.

- **Attributes + Features** We concatenate the predicted labels in Attribute-based model with CNN features as a combined representation.

- **KB - Affordance (Attributes)** A smaller KB learned without affordances (attributes).

- **KB - Full** Our full KB model defined in Sec. 3.2.

The Attributes + Features model (the third row in Table 3) outperforms the Attribute-based model (the second

row in Table 3) by 11.7%, indicating the importance of modeling continuous features in the KB. The full model KB - Full achieves the state-of-the-art performance on both basic-level and fine-grained classes with more than 2% improvement over the CNN baseline.

Fig. 6 offers some insight as to why a KB-based model performs well in a scene classification task. The class label is one of the many labels jointly inferred and predicted by the KB system, including attributes and affordances. So to predict an *auditorium*, attributes such as *indoor lighting*, *enclosed area*, and affordances such as *taking class for personal interest* can all help to reassure the prediction of an auditorium, and vice versa.

As mentioned in Sec. 3.2, we have collected annotations of 227 affordance classes for each of the 298 scene categories. We report the performance of *affordance prediction* by mean average precision (mAP) and mean F1 score (mF1) over the 227 affordance classes. The results are presented in Table 4. Here we compare our full KB model with the CNN Fine-tuned model [54], where we trained an $\ell_2$-logistic regression classifier on the CNN features for each of the 227 affordance classes. The KB - Full model outperforms the CNN baselines on both metrics.

Recall that the KB framework learns the weights of the relations between entities (e.g., scene classes, attributes and affordance, etc.) in a joint fashion. We can then examine the strength of these relations by looking at the factor weights of the underlying MRF. A large positive weight between two entities indicate a strong co-occurrence relation; whereas a large negative weight indicates a strong negative correlation. Fig. 7 provides examples of both the strongest and the weakest correlations between scene classes and attributes (Fig. 7(a)), as well as scene classes and affordances (Fig. 7(b)). For example, the KB has learned that the class *beach* has a strong co-occurrence relation with the attribute *sand*, and the class *railroad track* lacks correlation with the affordance *teaching*.

## 6.3. Image Search by Text Queries

Using the same model and framework, we can also query our KB for sets of images, instead of just one (Sec. 6.2), such as "*find me images of a sunny beach*." Here we use the same dataset as in Sec. 6.2. This task can also be expressed by a single query where the image is taken as variables (see the example in Sec. 5).

We randomly generate 100 queries of a single label (scene category, affordance or attribute), and 100 queries

| Class | auditorium | landing deck | candy store | basilica | swimming pool aquatic theater | bindery fly bridge |
|---|---|---|---|---|---|---|
| Affordances | community and social work, taking class for personal interest, religious practices, waiting, attending the performing arts | transportation and material moving work, in transit / traveling, military work | eating & drinking, food presentation, picking up / dropping off child, reading for personal interest, relaxing | eating & drinking, attending or hosting parties, volunteer work, community and social work, religious practices | entertainment / arts / design / sports / media work, personal care and service work, socializing | boating, watching fishing, tobacco use, executive work, farming / fishing and forestry work |
| Attributes | congregating, indoor lighting, spectating, enclosed area, glossy | transporting things or people, asphalt, natural light, far-away horizon, man-made | no horizon, cluttered space, dirty, eating, waiting in line | open area, natural light, sunny, man-made, vacationing | still water, diving, no horizon, natural light, congregating | metal, sunny, wire, man-made, natural light |

Figure 6: **Sample prediction results by the full KB model.** The ground-truth categories (in black) are shown in the first row. The first four images show examples of correct predictions from our KB model, and the last two show incorrect examples. As our model jointly infers multiple labels of an image, we show the predicted affordances (second row) in blue, and the predicted attributes (third row) in green.

| | | |
|---|---|---|
| 7.29 | beach | sand |
| 5.68 | creek | moist / damp |
| 5.65 | house | shingles |
| -3.29 | sun deck | flowers |
| -3.69 | apse indoor | vinyl / linoleum |
| -3.86 | gorge | man-made |

(a) Top weighted relations between categories and attributes

| | | |
|---|---|---|
| 13.8 | mountain snowy | hunting |
| 13.6 | mountain | participating in equestrian sports |
| 12.5 | orchard | physical care of children |
| -0.94 | call center | medical services |
| -0.95 | machine shop | collecting as a hobby |
| -1.04 | railroad track | teaching |

(b) Top weighted relations between categories and affordances

Figure 7: **Examples of the strongest and the weakest relations in the learned KB.** (a) Relations between scene classes (left column) and scene attributes (right column). (b) Relations between scene classes (left column) and scene affordances (right column). In both (a) and (b), the number at the beginning of each row indicates the actual factor weight in the underlying MRF. The more positive the number, the stronger the correlation. We show relations with the largest positive and negative weights in the KB. To be consistent with Fig. 6, we use the same color scheme for attributes and affordances.

of a pair of labels, each having at least 50 positive samples in the test set. Given a set of query labels, we aimed to retrieve the test images that are annotated with all the semantic labels in the set. We compare with two nearest neighbor baseline methods [18]. NNall ranks the test images based on the minimum Euclidean distance to any individual positive sample in the training set. NNmean ranks the images based on the distance to the centroids of the features of the positive samples. We report the mean precision at $k$, the mean fraction of correct retrievals out of the top $k$ over all queries, where $k$ goes from 1 to 50. As shown in Fig. 8, our method outperforms both simple nearest neighbor baselines when $k > 5$. NNmean performs better than ours among the top five retrievals; however, the false positive rate grows as the number of retrievals increases. In contrast, the relations
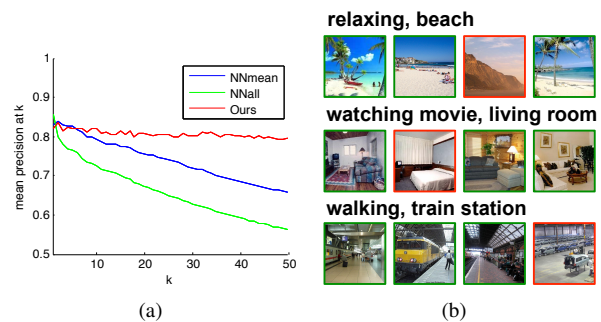


Figure 8: (a) **Performance variations of top $k$ retrievals** We compare our method with two nearest neighbor baselines. In contrast to these two methods, the KB model maintains a steady performance on lower-ranked retrievals. (b) **Top retrievals of example queries.** We show top four retrievals from three sample queries (in bold) by our KB model. The green boxes indicate correct retrievals, and red ones indicate incorrect retrievals.

in the KB compensate the weak and noisy visual signals, and, as a result, maintain stable and good performance on lower-ranked retrievals.

# 7. Conclusion

This paper presents a principled framework to perform learning and inference on a large-scale multimodal knowledge base (KB). Our contribution is to build a scalable KB to answer a variety of visual queries without re-training. Our KB is capable of making predictions on a number of standard vision tasks, on par with state-of-the-art models trained specifically for those tasks. In addition to these custom-trained classifiers, it is also interesting to explore these knowledge representations as an attempt towards tackling complex queries in real-world vision applications. Furthermore, this platform can be used to explore image-based reasoning. Towards these goals, future directions include a tighter integration between language and vision, and a more robust model for incorporating richer information.

# References

[1] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, L. Zitnick, and D. Parikh. VQA: Visual question answering. *ICCV*, 2015. 2

[2] J. Berant et al. Semantic parsing on Freebase from question-answer pairs. *EMNLP*, 2013. 5, 6

[3] K. Bollacker, C. Evans, P. Paritosh, T. Sturge, and J. Taylor. Freebase: A collaboratively created graph database for structuring human knowledge. *SIGMOD*, 2008. 1, 2, 13

[4] M. Á. Carreira-Perpiñán and G. E. Hinton. On contrastive divergence learning. *10th Int. Workshop on Artificial Intelligence and Statistics (AISTATS 2005)*, 2005. 3, 12

[5] Y.-W. Chao, Z. Wang, R. Mihalcea, and J. Deng. Mining semantic affordances of visual object categories. In *CVPR*, 2015. 2

[6] H. Chen and A. F. Murray. Continuous restricted boltzmann machine with an implementable training algorithm. In *Vision, Image and Signal Processing, IEE Proceedings*, volume 150, pages 153–158. IET, 2003. 3

[7] X. Chen, A. Shrivastava, and A. Gupta. NEIL: Extracting visual knowledge from web data. *ICCV*, 2013. 1, 2, 7

[8] X. Chen and C. L. Zitnick. Mind's eye: A recurrent visual representation for image caption generation. In *CVPR*, 2015. 2

[9] C. Desai, D. Ramanan, and C. Fowlkes. Discriminative models for multi-class object layout. In *ICCV*, 2009. 2

[10] S. Divvala, A. Farhadi, and C. Guestrin. Learning everything about anything: Webly-supervised visual concept learning. *CVPR*, 2014. 2, 7

[11] J. Donahue, L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell. Long-term recurrent convolutional networks for visual recognition and description. In *CVPR*, 2015. 2

[12] X. L. Dong et al. Knowledge Vault: A web-scale approach to probabilistic knowledge fusion. In *KDD*, 2014. 2, 4

[13] D. Ferrucci et al. Building Watson: An overview of the DeepQA project. *AI Magazine*, 2010. 1, 2, 4

[14] H. Gao, J. Mao, J. Zhou, Z. Huang, L. Wang, and W. Xu. Are you talking to a machine? dataset and methods for multilingual image question answering. *NIPS*, 2015. 2

[15] L. Getoor, N. Friedman, D. Koller, and A. Pfeffer. Learning probabilistic relational models. In *IJCAI*, 1999. 2, 4

[16] X. He and S. Gould. An exemplar-based crf for multi-instance object segmentation. In *CVPR*, 2014. 2

[17] X. He, R. S. Zemel, and M. Á. Carreira-Perpiñán. Multiscale conditional random fields for image labeling. In *CVPR*, volume 2, pages II–695. IEEE, 2004. 3

[18] K. Heller and Z. Ghahramani. A simple bayesian framework for content-based image retrieval. In *CVPR*, 2006. 8

[19] G. E. Hinton. Training products of experts by minimizing contrastive divergence. *Neural computation*, 14(8):1771–1800, 2002. 3, 12

[20] D. Hoiem, A. A. Efros, and M. Hebert. Putting objects in perspective. In *CVPR*, 2006. 2, 3

[21] A. Karpathy and L. Fei-Fei. Deep visual-semantic alignments for generating image descriptions. *CVPR*, 2015. 2

[22] D. Koller and N. Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009. 12

[23] C. Kong, D. Lin, M. Bansal, R. Urtasun, and S. Fidler. What are you talking about? text-to-image coreference. In *CVPR*, 2014. 2

[24] F. Kschischang, B. Frey, and H.-A. Loeliger. Factor graphs and the sum-product algorithm. *IEEE Trans. Inform. Theory*, 2001. 3, 12

[25] N. Kumar and S. Seitz. Photo recall: Using the internet to label your photos. In *VSM Workshop at CVPR*, 2014. 2

[26] S. Kumar and M. Hebert. Discriminative random fields: A discriminative framework for contextual interaction in classification. In *ICCV*, 2003. 2

[27] L. Ladicky, C. Russell, P. Kohli, and P. Torr. Graph cut based inference with co-occurrence statistics. *ECCV*, 2010. 2

[28] D. B. Lenat. Cyc: A large-scale investment in knowledge infrastructure. *Commun. ACM*, 1995. 2

[29] D. Lin, C. Kong, S. Fidler, and R. Urtasun. Generating multi-sentence lingual descriptions of indoor scenes. In *BMVC*, 2015. 2

[30] X. Lin and D. Parikh. Don't just listen, use your imagination: Leveraging visual common sense for non-visual tasks. *CVPR*, 2015. 2

[31] M. Malinowski and M. Fritz. A multi-world approach to question answering about real-world scenes based on uncertain input. In *NIPS*, 2014. 2

[32] M. Malinowski, M. Rohrbach, and M. Fritz. Ask your neurons: A neural-based approach to answering questions about images. *ICCV*, 2015. 2

[33] R. Mottaghi, X. Chen, X. Liu, N.-G. Cho, S.-W. Lee, S. Fidler, R. Urtasun, and A. Yuille. The role of context for object detection and semantic segmentation in the wild. *CVPR*, 2014. 2, 3

[34] F. Niu, C. Ré, A. Doan, and J. Shavlik. Tuffy: Scaling up statistical inference in Markov logic networks using an RDBMS. *Proc. VLDB Endow.*, 2011. 4

[35] F. Niu, B. Recht, C. Ré, and S. J. Wright. Hogwild!: A lock-free approach to parallelizing stochastic gradient descent. *NIPS*, 2011. 2, 5

[36] D. Parikh, C. L. Zitnick, and T. Chen. Exploring tiny images: The roles of appearance and contextual information for machine and human object recognition. *PAMI*, 2012. 2, 3

[37] G. Patterson, C. Xu, H. Su, and J. Hays. The SUN attribute database: Beyond categories for deeper scene understanding. *IJCV*, 2014. 4

[38] M. Richardson and P. Domingos. Markov logic networks. *Machine learning*, 2006. 1, 2, 3, 4, 5, 12

[39] F. Sadeghi, S. K. Divvala, and A. Farhadi. Viske: Visual knowledge extraction and question answering by visual verification of relation phrases. In *CVPR*, 2015. 2

[40] K. Shelley. Developing the american time use survey activity classification system. *Monthly Labor Review*, 2005. 4, 14

[41] B. Siddiquie, R. S. Feris, and L. S. Davis. Image ranking and retrieval based on multi-attribute queries. In *CVPR*, 2011. 2

[42] R. Socher, A. Karpathy, Q. V. Le, C. D. Manning, and A. Y. Ng. Grounded compositional semantics for finding and describing images with sentences. *TACL*, 2013. 2

[43] D. Suciu, D. Olteanu, C. Ré, and C. Koch. *Probabilistic Databases*. Morgan-Claypool, 2011. 5

[44] I. Sutskever, O. Vinyals, and Q. Le. Sequence to sequence learning with neural networks. *NIPS*, 2014. 2

[45] T. Tieleman. Training Restricted Boltzmann Machines using Approximations to the Likelihood Gradient. In *ICML*, 2008. 3

[46] A. Torralba, K. Murphy, and W. Freeman. Contextual models for object detection using boosted random fields. In *NIPS*, 2005. 2, 3

[47] K. Tu, M. Meng, M. W. Lee, T. E. Choe, and S.-C. Zhu. Joint video and text parsing for understanding events and answering queries. *IEEE MultiMedia*, 21(2):42–70, 2014. 2

[48] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and tell: A neural image caption generator. In *CVPR*, June 2015. 2

[49] M. Wick, A. McCallum, and G. Miklau. Scalable probabilistic databases with factor graphs and mcmc. *Proceedings of the VLDB Endowment*, 3(1-2):794–804, 2010. 3

[50] J. Xiao, J. Hays, K. Ehinger, A. Oliva, and A. Torralba. SUN database: Large-scale scene recognition from abbey to zoo. In *CVPR*, 2010. 2, 3, 4, 7, 14

[51] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhutdinov, R. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, 2015. 2

[52] B. Yao and L. Fei-Fei. Modeling mutual context of object and human pose in human-object interaction activities. In *CVPR*, 2010. 2

[53] L. Yu, E. Park, A. C. Berg, and T. L. Berg. Visual Madlibs: Fill in the blank Image Generation and Question Answering. *ICCV*, 2015. 2

[54] M. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *ECCV*, 2014. 3, 7, 11

[55] C. Zhang and C. Ré. DimmWitted: A study of main-memory statistical analytics. *Proc. VLDB Endow.*, 2014. 2, 5, 11

[56] Y. Zhao and S.-C. Zhu. Scene parsing by integrating function, geometry and appearance models. *CVPR*, 2013. 2

[57] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva. Learning Deep Features for Scene Recognition using Places Database. *NIPS*, 2014. 7

[58] J. Zhu, Z. Nie, X. Liu, B. Zhang, and J.-R. Wen. StatSnowball: a statistical approach to extracting entity relationships. In *WWW*, 2009. 1, 2

[59] Y. Zhu, A. Fathi, and L. Fei-Fei. Reasoning about object affordances in a knowledge base representation. *ECCV*, 2014. 1, 2, 3, 4, 5, 7

[60] C. L. Zitnick, D. Parikh, and L. Vanderwende. Learning the visual interpretation of sentences. In *ICCV*, 2013. 2

# A. Scalable Knowledge Base Construction

There are three key steps to make the knowledge base construction (KBC) scalable: data pre-processing, factor graph generation and high-performance learning. Sec. 4.1 provides an overview of the KBC process illustrating these three steps. Here we provide more detailed explanations of our knowledge base construction pipleine.

## A.1. Database Schema

The first step (the first box in Fig. 3) is to pre-process raw data into a structured representation. This representation enables us to perform structured queries (e.g. SQL) on the data. We provide the complete database schema in Fig. 9. The schema contains two types of tables: *data tables* contain the entities in Sec. 3.2 that are used to build the knowledge base (KB); *metadata tables* provide auxiliary information for the experiments and visualization. *sample_id* in Fig. 9 is a unique identifier of each training sample. These identifiers are used as a distribution key in the database system, where the data is distributed across segments as per the distribution keys.

Each data table stores entities of a certain type. We have a separate table for each of the four entity types in Sec. 3.2, where continuous values (image features) are stored as *double precision* numbers, and discrete values (scene category, affordance and attribute labels) are stored as *bigint*. We have seen in Sec. 4.1 that each row in the data tables corresponds to a variable in the factor graph. Thus the entities in Sec. 3.2 can be represented by different types of variables. We use 4096 continuous variables to represent an *Image* entity by its feature extracted from a fine-tuned CNN [54]. We use a multinomial variable to represent a scene category label, and Boolean variables to represent each of the attribute labels and affordance labels.

## A.2. Runtime environment

The knowledge base construction is conducted on a Non-Uniform Memory Access (NUMA) machine [55] with four NUMA nodes. Each has 12 physical cores and 24 logical cores, with Intel Xeon CPU@2.40GHz and 1TB main memory. We choose Greenplum as the underlying database system due to its power in massive parallel data processing.[4]

## A.3. Human-readable Rules

To define the KB with ease, we develop a declarative language, which serves as a human-readable interface for specifying the KB structure. The syntax of the declarative language is an extension to first-order logic in order to accommodate continuous variables. We introduced in Sec. 3.2 three types of relations. We define each type of relations by

---

[4] http://www.pivotal.io/big-data/pivotal-greenplum-database



| image features | | |
|---|---|---|
| id | *bigint* | |
| sample_id | *bigint* | * |
| dimension | *bigint* | |
| feat | *double precision* | |

| scene categories | | |
|---|---|---|
| id | *bigint* | |
| sample_id | *bigint* | * |
| category | *bigint* | |
| level | *bigint* | |

| scene affordances | | |
|---|---|---|
| id | *bigint* | |
| sample_id | *bigint* | * |
| affordance_id | *bigint* | |
| label | *bigint* | |

| scene attributes | | |
|---|---|---|
| id | *bigint* | |
| sample_id | *bigint* | * |
| attribute_id | *bigint* | |
| label | *bigint* | |

| scene attribute names | |
|---|---|
| attribute_id | *bigint* |
| name | *text* |

| train test split | |
|---|---|
| sample_id | *bigint* |
| holdout | *boolean* |

| scene affordance names | |
|---|---|
| affordance_id | *bigint* |
| name | *text* |

| scene categories names | |
|---|---|
| category | *bigint* |
| name | *text* |

Figure 9: **Database schema for structured representation.** The table names (in bold), column names (left) and data types (right) are provided. The blue boxes denote data tables containing KB entities; and the green ones denote metadata tables. The *id* column is a unique identifier for each row, which is used to create the factor graph. The stars (*) indicate the distribution keys for parallel data processing.

a group of rules, where each rule $R_j$ is a set specified with first-order logic formulas.

We first explain an example rule. We then describe the general form of the rules later. In Fig. 3 we have shown that our KBC system creates a factor in the factor graph of image I1 from the rule hasAffordance(I1, travel) ∧ hasAttribute(I1, sunny), which describes the co-occurrence between the affordance label travel and the attribute label sunny. We use the same example to show how factors are generated from the declarative language. Instead of writing rules for each of the affordance-attribute pair, we can simply write a rule:

$$\{(i, w(x,y), 1) \,|\, \text{hasAffordance}(i,x) \,\wedge\, \text{hasAttribute}(i,y)\}$$

where $i$, $x$ and $y$ correspond to the variables of images, affordance labels and attribute labels respectively. This rule can be instantiated by assigning values to these variables. One possible assignment is to set $i$ to image I1, $x$ to travel and $y$ to sunny. This creates a factor in the factor graph of image I1, where the factor value is 1 when hasAffordance(I1, travel) ∧ hasAttribute(I1, sunny) holds and 0 otherwise. It evaluates to 0 in the example of Fig. 3, as image I1 does not have attribute sunny. Under such variable assignment, the weight assigned to the factor is $w(\text{travel}, \text{sunny})$. It indicates that this weight will be shared by all the factors (one for each training image) that depict the co-occurrence between the affordance travel and the attribute sunny. This rule indicates that image I1

should have both `hasAffordance(I1,travel)` and `hasAttribute(I1,sunny)` to be true with a confidence score of $w(\text{travel}, \text{sunny})$. Similarly, the corresponding factors for other images share the same weight $w(\text{travel}, \text{sunny})$. More generally, each rule $R_j$ corresponds to a set in a given possible world $I$:

$$I(R_j) = \{(\bar{x}, w(\bar{y}), f(\bar{z}))\} \qquad (6)$$

where $\bar{x}, \bar{y}, \bar{z}$ are sets of variable in the domain (the set of all possible values the variables can take), and $w(\cdot)$ and $f(\cdot)$ are real-valued functions. Here $f(\cdot)$ essentially defines factors in the factor graph model and $w(\cdot)$ defines the corresponding factor weights (see Sec. 3.1). The arguments to $f(\cdot)$ define the variables required to compute the factor value. The arguments to $w(\cdot)$ define how the factor weights are shared across the factors.

All three types of relations in Sec. 3.2 can be specified as rules written in this declarative language. Fig. 10 provides a complete list of rules that we have used to build the visual KB. To be more specific, we express *image - label* relations using two sets of rules corresponding to 1) the linear terms, where the factors return the image feature values of each dimension; and 2) the bias terms, where the factors return a constant 1. For intra- and inter-correlations, we express them as conjunctions of two predicates, where the factors return 1 if both labels take the same Boolean value (either true or false), and 0 otherwise. In total, the proposed declarative language enables us to define the KB structure with eighteen first-order logic rules. Our KBC system automatically parses these rules, and creates a factor graph (see the second box in Fig. 3). Now we have the structure of the factor graph model, the next step is to learn the model parameters (i.e., factor weights). We will talk about the details of learning and inference in the next section.

## A.4. Learning and Inference

In this section we provide more technical details about learning and inference in our KB.

### A.4.1  Learning

The factor graph model in Sec. 3.1 is an instance of standard energy-based probabilisitic models [22] where the energy function $E(I)$ is defined through a linear combination of factors:

$$E(I) = \sum_{i=1}^{m} w_i f_i(I) \qquad (7)$$

A standard approach to learning is to optimize the negative log-likelihood of the training data in Eq. (3). Due to the intractability of computing the analytical gradients, sampling is a common practice to estimate the log-likelihood gradients. The gradient approximation used in Eq. (4) is a special case of contrastive divergence [19], called CD-1. Namely, instead of waiting for the Markov chain to converge, we obtain a sample after only one step of Gibbs sampling. This significantly reduces the cost of gradient computation per step, and has shown effective in several learning tasks [4, 19]. We illustrate in Fig. 3(d) that we create a factor graph for each image. This process is sometimes called *grounding* in the literature [38]. During training we treat these small factor graphs as a single large factor graph. The variables are mixed and shuffled before sampling. A weight update is performed at each Gibbs sampling step.

### A.4.2  Inference

The inference task is to derive the marginal probabilities of a conjunctive query in Eq. (5). This problem can be regarded as computing the expectation of a real function $f : \mathcal{I} \to \mathbb{R}$ given the probability distribution of possible worlds $I \in \mathcal{I}$:

$$\mathbf{E}[f; \mathbf{w}] = \sum_{I \in \mathcal{I}} \Pr[I; \mathbf{w}] f(I) \qquad (8)$$

where $\Pr[I; \mathbf{w}]$ is the probability of a possible world $I$ defined in Eq. (2), and $\mathcal{I}$ is the set of all possible worlds. Computing the exact expectation in Eq. (8) is intractable in general factor graphs, which requires summing over a large (or even infinite) number of variable assignments. Gibbs sampling is a commonly used method for approximate inference.

The Gibbs sampling starts with an initial world $I^{(0)}$. For each random variable $v_k$ in the factor graph, we sample its new value $v'_k$ from the conditional distribution $\Pr[v_k | MB(v_k); \mathbf{w}]$, where $MB(v)$ is the Markov blanket of the variable $v$. In the context of factor graphs [24], the Markov blanket of a variable is the set of factors that are connected to the variable. The sampler then moves to the next variable. After $m$ rounds of iterations, we have sampled a collection of possible worlds $\Omega = \{I^{(0)}, I^{(1)}, \dots, I^{(m)}\}$. We thus approximate the expectations of a query $q$ in Eq. (8) over $\Omega$:

$$\tilde{\mathbf{E}}[q] = \frac{1}{m} \sum_{i=1}^{m} q(I^{(i)}), \qquad (9)$$

where $q(I)$ is the value of the conjunctive query $q$ in possible world $I$. To be specific, $q(I)$ evaluates to 1 if all the predicates in the query $q$ are true in the possible world $I$, and 0 otherwise. After sufficient iterations, the probability of an answer to the query can be estimated by the number of iterations in which it takes that value over the total number of iterations.

## IMAGE−LABEL RELATIONS

**image features & scene category**
```
{(i,w(d),f) | sceneCategory(i,c) ∧ hasFeature(i,d,f)}
{(i,w(c),1) | sceneCategory(i,c)}
```

**image features & scene affordance**
```
scene_affordance_and_scene_features
{(i,w(a),f) | HasAffordance(i,a) ∧ hasFeature(i,d,f)}
{(i,w(a),1) | hasAffordance(i,a)}
```

**image features & scene attribute**
```
{(i,w(d),f) | hasAttribute(i,a) ∧ hasFeature(i,d,f)}
{(i,w(a),1) | hasAttribute(i,a)}
```

## INTRA−CORRELATIONS

**affordance & affordance**
```
{((i,a1,a2), w(a1,a2), 1) | hasAffordance(i,a1) ∧
    hasAffordance(i,a2)}
{((i,a1,a2), w(a1,a2), 1) | !hasAffordance(i, a1) ∧
    !hasAffordance(i, a2)}
```

**attribute & attribute**
```
{((i,a1,a2),w(a1,a2),1) | hasAttribute(i,a1) ∧
    hasAttribute(i,a2)}
{((i,a1,a2),w(a1,a2),1) | !hasAttribute(i,a1) ∧
    !hasAttribute(i,a2)}
```

## INTER−CORRELATIONS

**category & attribute**
```
{((i,c,a), w(a,c), 1) | sceneCategory(i, c) ∧
    hasAttribute(i, a)}
{((i,c,a), w(a,c), 1) | sceneCategory(i, c) ∧
    !hasAttribute(i, a)}
{((i,c,a), w(a,c), 1) | !sceneCategory(i, c) ∧
    hasAttribute(i, a)}
{((i,c,a), w(a,c), 1) | !sceneCategory(i, c) ∧
    !hasAttribute(i, a)}
```

**category & affordance**
```
{((i,c,a), w(a,c), 1) | sceneCategory(i, c) ∧
    hasAffordance(i, a)}
{((i,c,a), w(a,c), 1) | sceneCategory(i, c) ∧
    !hasAffordance(i, a)}
{((i,c,a), w(a,c), 1) | !sceneCategory(i, c) ∧
    hasAffordance(i, a)}
{((i,c,a), w(a,c), 1) | !sceneCategory(i, c) ∧
    !hasAffordance(i, a)}
```

Figure 10: **The complete list of rules for the visual knowledge base construction.** We build our visual knowledge base with the rules above. ! denotes negation and ∧ denotes conjunction. The formal semantics of the rules are described in Sec. A.3.

## B. Query Answering Application Setup

In Fig. 5, we have provided six query examples that illustrate the diversity of tasks our KB system can handle. In order to answer these diverse types of queries, it requires a fusion of information from various sources. In practice, we aggregate information from online databases, business and travel websites, etc. We provide the detailed experimental setups and the data sources here.

We augment our KB in Sec. 3.2 with a new set of geo-tagged images and several types of metadata. We briefly introduce the extra data sources that we used for this exper-

iment in Sec. 6.1. We randomly sample from Flickr100M[5] a pool of 20k images with geo-tags and timestamps. Besides these images, we incorporate additional information by either downloading from existing databases or crawling from the web. All the information is stored in a structured format as database tables (Sec. A.1).

1. We obtain a list of names and dates of 327 public holidays from Freebase[6] [3] from the instances of /time/holiday_category/holidays.

2. We scrape business information from Yelp.com and Hotels.com. We have crawled in total over sixteen thousand entries of business information, including 7k bars, 6k shopping centers and 3k hotels.

3. We download the daily temperature and weather data from National Climatic Data Center. Climate Data Online[7] (CDO) provides free access to global historical weather and climate data.

4. We download the publicly available GeoNames geographical database[8], which maps geolocations to over eight million place names.

We introduce new predicates in Fig. 11 (Boolean-valued functions) that enable us to query with these additional data. The semantics of these new predicates can be easily inferred from the predicate names and input variables. For instance, the predicate $\text{hasLocation}(img, \text{latlong1})$ evaluates to true if the image $img$ was annotated with the geo-location $\text{latlong1}$ and false otherwise; $\text{nearBy}(latlong1, latlong2, \text{1km})$ evaluates to true if the two geo-locations are within 1km away and false otherwise. Having defined the predicates, we use the augmented KB to answer the queries in Fig. 5. We list the conjunctive queries for each of the six example queries in Fig. 11. The predicates in each query are connected by logical conjunctions. Therefore the query evaluates to 1 if and only if every predicate in the query is true, and 0 otherwise. $\text{answer}(\cdot)$ indicates the return variables, i.e., the target answers to the queries. We retrieve a ranked list of the answers by computing a marginal probability of the queries (see Sec. 5 and Sec. A.4). Note that, once these additional metadata are incorporated into the KB framework, our system treats images, existing metadata and these new metadata on an equal footing in learning and inference. Therefore, a query can be answered by a joint inference with no post-filtering steps.

Following this approach, we are able to express richer and more complex queries by joining different pieces of information with logical conjunctions. As we can see, the

---

[5] http://yahoolabs.tumblr.com/post/89783581601/one-hundred-million-creative-commons-flickr-images
[6] https://www.freebase.com
[7] http://www.ncdc.noaa.gov/cdo-web/
[8] http://www.geonames.org/

**Q: Find me a modern looking mall near Fisherman's Wharf.**
hasLocation($img, latlong1$)
mall($mall, latlong2, zip$)
geoName(`Fisherman's Wharf`, $latlong3$)
hasAttribute($img$, `indoor lighting`)
hasAttribute($img$, `glossy`)
nearBy($latlong1, latlong2$, `1km`)
nearBy($latlong1, latlong3$, `20km`)
$\Rightarrow$ answer($img, mall, zip$)

**Q: Find me a place in Boston where I can play baseball.**
hasAffordance($img$, `playing baseball`)
hasLocation($img, latlong1$)
geoName(`Boston`, $latlong2$)
nearBy($latlong1, latlong2$, `1km`)
$\Rightarrow$ answer($img, latlong1$)

**Q: Find me a hotel in Boston with new furniture.**
hasLocation($img, latlong1$)
hasAttribute($img$, `glossy`)
geoName(`Boston`, $latlong2$)
nearBy($latlong1, latlong2$, `20km`)
hotel($hotel, latlong2, date, price, phone$)
$\Rightarrow$ answer($img, hotel, price, phone$)

**Q: Find me a cozy bar to drink beer near the AT&T Plaza.**
hasAttribute($img$, `cluttered space`)
hasLocation($img, latlong1$)
bar($bar, latlong2, price, phone$)
geoName(`AT&T Plaza`, $latlong3$)
nearBy($latlong1, latlong2$, `1km`)
nearBy($latlong1, latlong3$, `1km`)
$\Rightarrow$ answer($img, bar, price, phone$)

**Q: Find me a sunny and warm beach during Christmas Day 2013.**
sceneCategory($img$, `beach`)
hasAttribute($img$, `sunny`)
hasAttribute($img$, `warm`)
hasLocation($img, latlong1$)
geoName($location, latlong2$)
nearBy($latlong1, latlong2$, `1km`)
temperature($location, degree$, `2013/12/25`)
$\Rightarrow$ answer($img, location, degree, latlong2$)

**Q: Find me pictures of sunny days of Seattle during August.**
hasAttribute($img$, `sunny`)
hasLocation($img, latlong1$)
hasDate($img, day$, `August`, $year$)
geoName(`Seattle`, $latlong2$)
nearBy($latlong1, latlong2$, `20km`)
$\Rightarrow$ answer($img, day$, `August`, $year$)

Figure 11: Conjunctive queries for the query answering examples in Fig. 5. We omit the conjunction symbols ($\wedge$) between predicates for neatness.
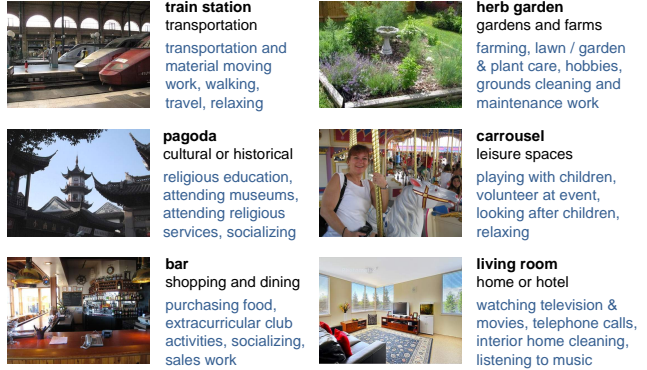


Figure 12: **Sample affordance annotations in the augmented scene dataset.** We augment the SUN dataset [50] with a lexicon of 227 affordances. We provide the fine-grained category (in bold), the basic-level category and a subset of their affordance annotations.

query language in Sec. 5 is capable of expressing a wide range of queries. Moreover, these queries can be answered in a principled manner, by evaluating marginals in the joint probability model. Given such a flexible framework, data becomes the key to extend our model's power of answering real-world questions. We are interested in exploring more efficient and automatic ways to aggregate information from large-scale multimodal corpora for future work.

## C. Affordance Annotations

We augment the SUN dataset [50] with additional annotations of scene affordances. We use a lexicon of 227 affordances (actions) from the American Time Use Survey (ATUS) [40] sponsored by the Bureau of Labor Statistics, which catalogs the actions in daily lives and represents United States census data. The original ATUS lexicon includes 428 specific activities organized into 17 major activity categories and 105 mid-level categories. We re-organize the categories by collapsing visually similar superordinate categories into one action. For instance, the superordinate-level category "traveling" was collapsed into a single category because being in transit to go to school should be visually indistinguishable from being in transit to go to the doctor. This results in 227 actions in total. Fig. 12 shows six example images with a subset of their affordance annotations.

The lexicon covers a broad space of possible actions that could take place in scenes. We conducted a large-scale online experiment with over 400 AMT workers annotating the possibilities of the 227 actions for each of the 298 scene categories (Sec. 3.2). 10 votes are collected for each category-affordance pair. Positive ($\geq$ 3 votes) and negative ($\leq$ 2 votes) annotations are selected as evidence. These 227 affordances are listed in alphabetic order below:

**A** appliance repair & maintenance (self), architecture and engineering work, arts & crafts, arts & crafts with children, arts / design / entertainment / sports / media work, attending child's events, attending meetings for personal interest, attending movies, attending museums, attending or hosting parties, attending religious services, attending school-related meetings & conferences, attending the performing arts

**B** banking, biking, boating, bowling, building & repairing furniture, building and grounds cleaning and maintenance work, business and financial operations work, buying / selling real estate

**C** camping, civic obligations, cleaning home exterior, collecting as a hobby, community and social work, comparison shopping, computer and mathematical work, computer use (not games), construction and extraction work

**D** dancing, doing aerobics, doing gymnastics, doing martial arts

**E** eating & drinking, education and library work, education-related administrative activities, email, exercising & playing with animals, exterior home repair & decoration, extracurricular club activities

**F** farming / fishing and forestry work, fencing, financial management, fishing, food & drink preparation, food preparation and serving work, food presentation

**G** gambling, golfing, grocery shopping

**H** health-related self care, healthcare work, helping adult, helping child with homework, hiking, hobbies, home heating / cooling, home security, home-schooling children, homework, household organization & planning, hunting

**I** in transit / traveling, income-generating hobbies & crafts, income-generating performance, income-generating rental property activity, income-generating selling activities, income-generating services, installation / maintenance and repair work, interior decoration & repair, interior home cleaning

**J** job interviewing, job search activities

**K** kitchen & food clean-up

**L** laundry, lawn / garden & plant care, legal work, listening to music (not radio), listening to radio, looking after adult, looking after children

**M** mailing, maintaining home pool / pond / hot tub, management / executive work, military work

**N** non-veterinary pet care

**O** obtaining licenses & paying fees, obtaining medical care for adult, obtaining medical care for child, office and administrative work, organizing & planning for adults, organizing & planning for children, out-of-home medical services

**P** participating in aquatic sports, participating in equestrian sports, participating in rodeo, personal care and service work, physical care of adults, physical care of children, picking up / dropping off adult, picking up / dropping off child, playing baseball, playing basketball, playing billiards, playing football, playing games, playing hockey, playing racquet sports, playing rugby, playing soccer, playing softball, playing sports with children, playing volleyball, playing with children (not sports), production work, protective services work, providing medical care to adult, providing medical care to child, purchasing food (not groceries), purchasing gasoline

**R** reading for personal interest, reading with children, relaxing, religious education, religious practices, rock climbing / caving, rollerblading / skateboarding, running

**S** sales work, school music activities, science work, security screening, sewing & repairing textiles, sexual activity, shopping (except food and gas), skiing / ice skating / snowboarding, sleeping, socializing, storing household items, student government

**T** taking class for degree or certification, taking class for personal interest, talking with children, telephone calls, tobacco use, transportation and material moving work, travel, using cardiovascular equipment

**U** using clothing repair & cleaning services, using home repair & construction services, using in-home medical services, using interior home cleaning services, using lawn & garden services, using legal services, using meal preparation services, using other financial services, using paid childcare services, using personal care services, using pet services, using police & fire services, using professional