

Retrieving Black-box Optimal Images from External Databases

Ryoma Sato

r.sato@ml.ist.i.kyoto-u.ac.jp
Kyoto University / RIKEN AIP

ABSTRACT

Suppose we have a black-box function (e.g., deep neural network) that takes an image as input and outputs a value that indicates preference. How can we retrieve optimal images with respect to this function from an external database on the Internet? Standard retrieval problems in the literature (e.g., item recommendations) assume that an algorithm has full access to the set of items. In other words, such algorithms are designed for service providers. In this paper, we consider the retrieval problem under different assumptions. Specifically, we consider how users with limited access to an image database can retrieve images using their own black-box functions. This formulation enables a flexible and finer-grained image search defined by each user. We assume the user can access the database through a search query with tight API limits. Therefore, a user needs to efficiently retrieve optimal images in terms of the number of queries. We propose an efficient retrieval algorithm TIARA for this problem. In the experiments, we confirm that our proposed method performs better than several baselines under various settings.

CCS CONCEPTS

• **Information systems** → **Web searching and information discovery**; **Users and interactive retrieval**.

KEYWORDS

Information Retrieval, Web Searching, Linear Bandits, Private Recommender Systems

ACM Reference Format:

Ryoma Sato. 2022. Retrieving Black-box Optimal Images from External Databases. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining (WSDM '22)*, February 21–25, 2022, Tempe, AZ, USA. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3488560.3498462>

1 INTRODUCTION

As the amount of information on the Internet continues to drastically increase, information retrieval algorithms are playing more important roles. In a typical situation, a user of a system issues a query by specifying keywords, and an information retrieval algorithm retrieves the optimal items with respect to the query words. Here, the retrieval algorithm is designed by the service provider, not by the user. The uses of information systems have become

divergent, and various retrieval algorithms have therefore been proposed, e.g., a cross-modal image search [9, 32], complex query retrieval [47], and conversational recommender systems [16, 58].

Deep neural networks have achieved state-of-the-art performances in computer vision tasks [27, 34], notably image retrieval [4, 8, 25, 48]. In a conventional setting of image retrieval, algorithms assume that they have full access to the image database. A straightforward method under this setting is to evaluate all images in the database and return the one with the maximum score. When the image database is extremely large, two-step methods are used, i.e., a handful of images are retrieved through fast retrieval algorithms such as a nearest neighbor search, and the results are ranked using sophisticated algorithms. Hash coding further improves the effectiveness and efficiency [36, 40].

In this study, we consider information retrieval under a different scenario. Whereas most existing studies have focused on how a service provider can improve the search algorithms, we focus on how a user of a service can effectively exploit the search results. Specifically, we consider a user of a service builds their own scoring function. The examples of the scoring functions are as follows.

Example 1 - Favorite Image Retrieval: A user trains deep neural networks using a collection of favorite images found in different services and wants to retrieve images with similar properties in a new service.

Example 2 - Similar Image Retrieval: Deep convolutional neural networks are known to have a superior ability to extract useful image features [4, 8, 23]. Although some online services provide similar image search engines, users do not have full control of the search. For example, even if the service provides a similar texture-based image search engine, some users may want to retrieve similar images based on the semantics. The user-defined score function allows image searches at a finer granularity.

Example 3 - Fair Image Retrieval: A search engine can be unfair to some protected attributes. For example, when we search for images through a query “president,” an image search engine may retrieve only male president images [56]. Some users may want to use their own scoring function, e.g., scoring male and female images equally.

Recent advancements in deep learning, such as self-supervised contrastive learning [13, 28] and meta-learning [22], enable the training of deep neural networks with a few labeled samples. In addition, many pre-trained models for various vision tasks have been released on the Internet. Machine learning as a service platform, such as Microsoft Azure Machine Learning Studio and AWS Machine Learning, has also made it easy to build a machine learning model. These techniques and services enable an individual to easily build their own black-box scoring function. Suppose a user has already built a black-box function she wants to optimize. How can she retrieve optimal images with respect to this function from an image database on the Internet?

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

WSDM '22, February 21–25, 2022, Tempe, AZ, USA

© 2022 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-9132-0/22/02.

<https://doi.org/10.1145/3488560.3498462>

Table 1: Symbols, Definitions, and Examples.

| Symbol | Definition | Example |
|---|---|--|
| \mathcal{X} | The set of images in the image database | The set of all images uploaded on Flickr |
| \mathcal{T} | The set of tags | The set of all tags on Flickr |
| $f: \mathcal{X} \rightarrow \mathbb{R}$ | The objective black-box function | A deep neural network that computes a preference score of an image |
| $O: \mathcal{T} \rightarrow \mathcal{X} \times 2^{\mathcal{T}}$ | The oracle for image search | A wrapper of flickr.photos.search API |
| $\mathcal{T}_{\text{ini}} \subseteq \mathcal{T}$ | Initially known tags | 100 popular tags obtained by flickr.tags.getHotList API |
| $B \in \mathbb{Z}_+$ | A budget of oracle calls | $B = 500$ |

Under this setting, an individual cannot evaluate all images in the image database because it contains a significant number of images. It is also impossible for an individual to build a search index (e.g., a hash index) because of both limited access to the database and the insufficient computational resources of the individual. These limitations prohibit the use of standard image retrieval algorithms.

In this paper, we assume that a user can access the database through a search query alone and that a tight query budget exists. In addition, we assume little knowledge about the database to apply our method to a new environment. We formulate the problem through the lens of the multi-armed bandit problem and propose a query efficient algorithm, TIARA, with the aid of pre-trained representative word embeddings. We confirm the effectiveness of the proposed method under various settings, including the online Flickr environment.

It is noteworthy that Bayesian optimization and optimization on deep neural networks [21, 46, 55] also aim to optimize black-box or complex functions. However, these methods assume a continuous and simple optimization domain, typically the entire Euclidean space \mathbb{R}^d , a hypercube $[0, 1]^d$, or a unit ball $\{x \mid \|x\| \leq 1\}$. By contrast, we aim to retrieve optimal images from the *fixed database*. Therefore, the optimization domain is a discrete set of images. Although there are several methods for discrete black-box optimization [6, 49], they also assume that the optimization domain is simple, e.g., $\{0, 1\}^d$, or at least they have full access to the optimization domain. Under our setting, an algorithm does not even know the entire optimization domain, i.e., the image database. This limitation causes numerous challenges, as we describe in the following sections.

The contributions of this paper are summarized as follows:

- We formulate the black-box optimal image retrieval problem. This problem examines how a user of an online service can effectively exploit a search engine, whereas most existing studies focus on how a service provider can improve a search engine.
- We propose TIARA, an effective method for this problem for the first time. TIARA is a general algorithm that works in various situations and retrieves optimal black-box images with few API queries.
- We investigate the effectiveness of TIARA using many real-world data. Notably, we conduct “in the wild” experiments on the real-world Flickr environment and confirm that TIARA can be readily used in real-world applications.

Reproducibility: Our code is available at <https://github.com/joisino/tiara>.

2 BACKGROUND

2.1 Notations

Let \mathcal{X} be a set of images in the image database. We do not know the exact set of \mathcal{X} . Let $f: \mathcal{X} \rightarrow \mathbb{R}$ be a black-box function that evaluates the value of an image. For example, f measures the preference of the user in favorite image retrieval and measures the similarity in similar image retrieval. The notations are summarized in Table 1.

2.2 Problem Formulation

In this section, we formulate the problem of black-box optimal image retrieval. We observe that many image databases, such as Flickr and Instagram, support (hash) tag-based search. We assume that we can search for images by specifying a tag through an API. For example, in the Flickr case, we use the flickr.photos.search API.

For a formal discussion, we generalize the function of tag search APIs. Let $\mathcal{T} \subset \Sigma^*$ be the set of tags, where Σ^* is the set of strings. We formalize a tag search API as a (randomized) oracle O that takes a tag as input and returns an image with the tag. We assume that O always returns different images even when we query the same tag twice. In addition to the image itself, we assume that the oracle O returns the set of tags of the image. Therefore, for any tag $t \in \mathcal{T}$, $O(t) \in \mathcal{X} \times 2^{\mathcal{T}}$. Let $O(t).\text{image} \in \mathcal{X}$ denote the returned image and $O(t).\text{tags} \in 2^{\mathcal{T}}$ denote the returned tags. As the oracle returns an image with the query tag, $t \in O(t).\text{tags}$ always holds.

We find that myriad tags exist in real-world services, and we cannot know the entire tag set. Therefore, we assume that we know only a fraction \mathcal{T}_{ini} of the tag set. Here, \mathcal{T}_{ini} can be constructed by browsing the online service or retrieving popular tags through an API, e.g., flickr.tags.getHotList, in the Flickr example.

We assume that there is a budget $B \in \mathbb{Z}_+$ for the oracle call. For example, in the Flickr case, There is an API rate limit of 3600 calls per hour. Thus, it is natural to assume that we can use at most 3600 API calls in one task. If we retrieve many images within a short period of time, the API limits become tighter for each retrieval. To summarize, the black-box optimal image retrieval problem is formalized as follows.

Black-box optimal image retrieval.

Given:

- Black box function $f: \mathcal{X} \rightarrow \mathbb{R}$
- Oracle $O: \mathcal{T} \rightarrow \mathcal{X} \times 2^{\mathcal{T}}$
- Known tags $\mathcal{T}_{\text{ini}} \subseteq \mathcal{T}$
- Budget $B \in \mathbb{Z}_+$

Goal: Find an image $x \in \mathcal{X}$ in the image database with as high an $f(x)$ as possible within B accesses to the oracle.

3 PROPOSED METHOD

In this section, We introduce our proposed method, TIARA (Tag-based image retrieval).

3.1 Bandit Formulation

We first propose regarding the black-box optimal image retrieval problem as a multi-armed bandit problem [37, 57]. Specifically, we regard a tag as an arm, the budget B as the time horizon, and the objective value $f(O(t).\text{image})$ as the reward when we choose arm t . This formulation enables us to use off-the-shelf multi-armed bandit algorithms, such as UCB [3], ϵ -greedy, and Thompson sampling [60]. This formulation is the basics of our proposed algorithm. However, there are several challenges to a black-box optimal image retrieval problem. First, we do not initially know all arms but only a fraction \mathcal{T}_{ini} of arms. Considering \mathcal{T}_{ini} alone leads to suboptimal results because \mathcal{T}_{ini} does not contain the best arm in general. Second, myriad arms exist, and the number $|\mathcal{T}|$ of arms is larger than budget B in practice. Standard multi-armed bandit algorithms first explore all arms once. However, they are unable to even finish this initial exploration phase under a tight budget constraint.

The first problem is relatively easy to solve. We obtain tags $O(t).\text{tags}$ when we choose arm t . These tags may contain new tags. We can add such tags into the known tag set and gradually grow the known set. A good bandit algorithm will choose relevant tags, and the returned tags will contain relevant tags. For example, suppose that the black-box function f prefers cat images. A good bandit algorithm will choose “cat” and “animal” tags and obtain cat images accompanied by many cat-related tags. Even if the budget is so tight that we cannot know all tags within the budget, a good bandit algorithm ignores irrelevant tags and prioritizes the collection of many relevant tags. The second problem is essential and difficult to solve. We investigate how to improve the query efficiency using tag embeddings in the following sections.

Discussion on the max-bandit problem. It should be noted that standard bandit algorithms aim to maximize a cumulative reward or regret, whereas the black-box optimal image retrieval problem aims to find an image with the maximum value. These goals are not exactly the same. The problem of obtaining a maximum reward is known as a max K -armed bandit [17] or extreme bandit [10]. However, we adhere to the standard bandit setting in this study for the following reasons: First, we assume that the reward distribution is Gaussian-like, whereas existing max- K -armed bandit algorithms [10, 17] mainly assume that rewards are drawn from extreme value distributions, such as the GEV distribution. We find that this is not the case in the applications considered herein, i.e., image retrieval. Second, under our setting, an arm with a high expected value often leads to a high maximum value because such an arm usually represents the concept captured by the black-box function f . By contrast, max-bandit algorithms focus on detecting heavy tail arms with possibly low average rewards. We find that this is not necessary for our setting. Finally, we experimentally confirm that our proposed method already shows a superior empirical performance even though it does not directly maximize the maximum objective value. We leave the exploration of the max-bandit algorithms in our setting for future work.

Algorithm 1: TIARA

```

1 Data: Known tags  $\mathcal{T}_{\text{ini}} \subseteq \mathcal{T}$ , Black box function  $f: \mathcal{X} \rightarrow \mathbb{R}$ ,
   Budget  $B \in \mathbb{Z}_+$ , Oracle  $O: \mathcal{T} \rightarrow \mathcal{X} \times 2^{\mathcal{T}}$ , Tag
   Embedding  $\{v_t \in \mathbb{R}^d \mid t \in \mathcal{T}\}$ , Regularization
   coefficient  $\lambda \in \mathbb{R}_+$ , Exploration coefficient  $\alpha \in \mathbb{R}_+$ .
2 Result: An image  $x \in \mathcal{X}$  in the image database with as high
    $f(x)$  as possible.
3  $A \leftarrow \lambda I_{d \times d}$  // Initialize  $A$ 
4  $b \leftarrow 0_d$  // Initialize  $b$ 
5  $\mathcal{T}_{\text{known}} \leftarrow \mathcal{T}_{\text{ini}}$  // Initialize with the initial tags
6 for  $i \leftarrow 1$  to  $B$  do
7    $s_t \leftarrow v_t^\top A^{-1} b + \alpha \sqrt{v_t^\top A^{-1} v_t} \quad \forall t \in \mathcal{T}_{\text{known}}$  // Scores
8    $t_i \leftarrow \operatorname{argmax}_{t \in \mathcal{T}_{\text{known}}} s_t$  // Choose the best arm
9    $r_i \leftarrow O(t_i)$  // Issue a query
10  for  $t \in r_i.\text{tags}$  do
11     $A \leftarrow A + v_t v_t^\top$  // Update  $A$ 
12     $b \leftarrow b + f(r_i.\text{image}) \cdot v_t$  // Update  $b$ 
13     $\mathcal{T}_{\text{known}} \leftarrow \mathcal{T}_{\text{known}} \cup \{t\}$  // Insert  $t$  to  $\mathcal{T}_{\text{known}}$ 
14 return  $\operatorname{argmax}_{x \in \{r_i.\text{image} \mid i=1, \dots, B\}} f(x)$  // Best image

```

3.2 Tag Embedding

Because we cannot choose each arm even once on average, we need to estimate the value of each arm without observing the reward from it. We estimate the value of one arm from the rewards obtained from the other arms. The key is how to define the similarities of the arms. As the challenge here, we assume a new image database, and owing to the tight API limit, it is difficult to learn the similarities from the current environment in an online manner. To tackle this problem, we utilize external resources. We find that a tag in an image database is usually described through natural language and typically is a word or composition of words. We use a pre-trained word embedding [45, 50] to define the similarities between tags. Specifically, we first decompose a tag into a bag of words by non-alphabetical symbols, e.g., a white space. We use the mean of the word embeddings in the bag of words as the tag embedding.

Owing to the tag embeddings, we can infer the value of an arm from similar arms. For example, if the reward from a “cat” tag is high (resp. low), we can assume the black-box function is highly (resp. rarely) related to cats, and we can infer that similar tags with similar embeddings, such as “Aegean cat” and “animal,” are also valuable (resp. irrelevant) without actually querying them.

3.3 Tiara

Our proposed method combines the aforementioned tag embeddings with a bandit algorithm. Specifically, we utilize LinUCB [39]. Because there are no contextual features, we use only features of arms. We define the feature of an arm as the tag embedding introduced in Section 3.2 and apply LinUCB to this feature. We call this variant TIARA-S, where S stands for “single” and “simple.” However, we found that there are too many tags, and LinUCB is still inefficient for learning a reward because of the tight query budget and limited training samples. To improve the query efficiency, we use another signal from the oracle. The oracle O returns not only

an image but also the tags of the returned image. We assume that these tags have similar average rewards and add these tags into the training dataset. Specifically, when we query tag $t \in \mathcal{T}$, we use $\{(s, f(O(t).\text{image})) \mid s \in O(t).\text{tags}\}$ as training data, whereas TIARA-S uses only $\{(t, f(O(t).\text{image}))\}$. As we will see in the experiment, this technique significantly improves query efficiency and performance. Algorithm 1 shows the pseudo-code of TIARA.

Time Complexity. The bottleneck of the computation is in Line 7. The inverse matrix A^{-1} can be efficiently computed in an iterative manner using the Sherman–Morrison formula, i.e.,

$$(A + v_t v_t^\top)^{-1} = A^{-1} - \frac{A^{-1} v_t v_t^\top A^{-1}}{1 + v_t^\top A^{-1} v_t}.$$

This technique reduces the cubic dependence on the number of dimensions to the quadratic dependence. Therefore, the computation is in Line 7 takes $O(|\mathcal{T}_{\text{known}}|d^2)$ time. Let T_{max} be the maximum number of tags of an image. The value of T_{max} is typically 10 to 100. Because $|\mathcal{T}_{\text{known}}|$ increases by at most T_{max} in an iteration, the total time complexity is $O((|\mathcal{T}_{\text{ini}}| + BT_{\text{max}})Bd^2)$ in the worst case. In our problem setting, B is small (e.g., in the hundreds) because of the tight API limitation, and $|\mathcal{T}_{\text{ini}}| \approx 100$ and $d \approx 300$ are also small during the experiments. Therefore, the oracle calls in Line 9 become a bottleneck in the wall-clock time because it requires communication to the Internet. As the oracle calls are common in all methods, TIARA is sufficiently efficient. When the black-box function is complex, its evaluation can be a computational bottleneck as well. TIARA evaluates f as few as B times, which is the same as with other baseline methods. In addition, when the efficiency is insufficient, we can speed up TIARA through a lazy variance update [19] and by applying sub-sampling heuristics.¹

Discussion on graph-feedback bandits. There are several multi-armed bandit algorithms with a so-called feedback graph setting [1, 31, 42]. This is an intermediate setting between the bandit and full feedback. Specifically, it assumes that there is an underlying graph, where a node represents an arm, and when we choose arm t , we can also observe the rewards of the neighboring arms. The underlying graph can be directed and time-varying. Therefore, if we define the neighbor of t as $O(t).\text{tags}$, our problem is seen as a variant of the feedback graph setting. However, we do not employ the feedback graph framework for the following reasons: First, existing graph feedback bandit algorithms require the number of neighbors of each neighbor [1, 31, 42], i.e., $|f(O(s).\text{tags})|, \forall s \in f(O(t).\text{tags})$ in our setting. We need additional API queries to compute these values. Such additional queries are prohibitive because the API limit is tight under our setting. Second, graph feedback bandit algorithms assume that the reward feedback of neighboring arms is an unbiased estimate of the true rewards [1, 31, 42]. However, $f(O(t).\text{image})$ is not unbiased for arm $s \in O(t).\text{tags}$. Therefore, we cannot enjoy the theoretical guarantees of the graph feedback bandits. Third, the improvement obtained by the feedback graph framework is on the order of $O(\sqrt{\alpha/n})$ [31], where n is the number of arms, and α is the size of the maximum independent set of the feedback graph. Under our setting, the sizes of the feedback graphs and α are still large

in practice. How to leverage the existing graph-feedback bandit algorithms for our setting is an interesting area of future study.

Discussion on the reward estimation model. TIARA uses LinUCB and a linear model for estimating the reward from the feature vector. In general, we can use any model for prediction. We use a linear model owing to the following reasons: First, complex models are difficult to train with a tight sample budget. Second, LinUCB empirically performs quite well under various tasks [33, 39] regardless of its suboptimal theoretical regret. Third, we use the average word embeddings as the feature vector of an arm. It has been confirmed that word embeddings are representative and that the weighted average of the word embeddings with a cosine similarity or linear models produce superior performances in many NLP tasks [2, 54] and sometimes outperform even neural network-based methods.

3.4 Visualization and Interpretation

Although not our original goal, TIARA provides an interpretation of the black-box function f as a byproduct.

Deep neural networks suffer from interpretability issues for reliable decision making. There have been many interpretation methods developed for deep neural networks [55, 59]. We consider model-level interpretability [46, 55, 62], i.e., interpreting what function each model represents. It is unclear what function each model represents by simply looking at the model parameters of deep models. Even if the model is prepared by the user, deep models sometimes behave in unexpected ways [24, 41]. Input instances that produce high values can be regarded as representations of the model [46, 55, 62]. TIARA can retrieve such images from external image databases. Compared to methods that use a fixed dataset, the search space of TIARA is extremely large. Therefore, there are more chances that relevant images will be found. In addition, tag scores of TIARA also provide another interpretation of the black-box function through words. We show that these tags are also beneficial for interpretability by visualizing word clouds in the experiments. The tag-based interpretability is beneficial for further exploration as well. When the result is unsatisfactory, the user can continue manual exploration from the tags with high scores.

4 EXPERIMENTS

We investigate the performance of TIARA using various real-world datasets.

4.1 Experimental Setups

4.1.1 Baselines. We use the following baselines.

- **Random** queries random known tags.
- **ϵ -Greedy** is a bandit algorithm. This algorithm chooses the best tag with the highest empirical mean reward with a probability of $1 - \epsilon$ and chooses a random tag with probability ϵ . The candidate pool is set to \mathcal{T}_{ini} .
- **UCB** is a bandit algorithm. The score of tag t is the empirical mean reward plus $\alpha\sqrt{1/n_t}$, where n_t is the number of observations from tag t , and α is a hyperparameter. UCB chooses the tag with the highest score. The candidate pool is set to \mathcal{T}_{ini} .

¹We do not adopt such techniques in the experiments.



Figure 1: Open Image Environment. (Top) Tench, (Bottom) Black Swan. (Left) Learning curves averaged over 10 independent runs. (Mid) Word clouds that represent the scores computed by TIARA at the end of the last iteration. These visualizations provide an interpretability of the black-box function. (Right) Best images found at each iteration.

Table 2: Open Image Environment. Each score is the highest $f(x)$ found by an algorithm. The scores are averaged over 10 independent runs, and the standard deviations are also reported. The highest score is highlighted by bold in each column.

| Method \ Class | Tench | Black Swan | Tibetan Terrier | Tiger Beetle | Academic Gown | Cliff Dwelling | Hook | Paper Towel | Slot Machine | Water Tower | Average |
|--|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|
| Random | 13.11 ± 5.95 | 12.62 ± 4.83 | 13.16 ± 1.92 | 8.21 ± 1.85 | 12.63 ± 4.88 | 12.63 ± 5.31 | 11.58 ± 1.51 | 10.40 ± 1.01 | 9.29 ± 1.68 | 12.09 ± 1.64 | 11.57 ± 0.89 |
| ϵ -Greedy ($\epsilon = 0.01$) | 13.48 ± 6.27 | 12.30 ± 4.94 | 13.76 ± 2.16 | 7.47 ± 2.53 | 14.83 ± 5.80 | 13.32 ± 6.07 | 11.17 ± 1.85 | 9.79 ± 0.84 | 8.55 ± 1.02 | 13.25 ± 1.41 | 11.79 ± 0.99 |
| ϵ -Greedy ($\epsilon = 0.1$) | 13.48 ± 6.27 | 12.06 ± 4.99 | 13.76 ± 2.16 | 8.54 ± 3.19 | 14.83 ± 5.80 | 13.06 ± 6.26 | 11.17 ± 1.85 | 9.66 ± 0.83 | 8.84 ± 0.88 | 13.25 ± 1.41 | 11.86 ± 1.08 |
| ϵ -Greedy ($\epsilon = 0.5$) | 13.64 ± 6.07 | 12.30 ± 4.96 | 13.76 ± 2.16 | 9.88 ± 3.39 | 13.10 ± 4.77 | 13.36 ± 6.02 | 11.64 ± 1.96 | 9.74 ± 0.95 | 8.42 ± 1.16 | 13.06 ± 1.63 | 11.89 ± 1.07 |
| UCB ($\alpha = 0.1$) | 13.48 ± 6.27 | 12.30 ± 4.94 | 13.76 ± 2.16 | 7.47 ± 2.53 | 14.83 ± 5.80 | 13.32 ± 6.07 | 11.17 ± 1.85 | 9.79 ± 0.84 | 8.55 ± 1.02 | 13.25 ± 1.41 | 11.79 ± 0.99 |
| UCB ($\alpha = 1.0$) | 13.71 ± 5.99 | 12.33 ± 4.94 | 13.76 ± 2.16 | 9.27 ± 3.73 | 14.83 ± 5.80 | 13.32 ± 6.07 | 11.17 ± 1.85 | 10.14 ± 0.98 | 8.42 ± 1.17 | 13.25 ± 1.41 | 12.02 ± 1.14 |
| UCB ($\alpha = 10.0$) | 13.68 ± 5.98 | 12.62 ± 4.83 | 13.76 ± 2.16 | 8.21 ± 1.85 | 13.05 ± 4.58 | 13.36 ± 6.02 | 11.61 ± 2.00 | 9.77 ± 1.25 | 8.70 ± 1.50 | 13.25 ± 1.41 | 11.80 ± 0.87 |
| Ada- ϵ -Greedy ($\epsilon = 0.1$) | 11.34 ± 5.84 | 12.81 ± 4.97 | 13.28 ± 1.97 | 8.70 ± 1.36 | 15.62 ± 3.38 | 10.74 ± 2.13 | 11.09 ± 1.72 | 9.93 ± 0.81 | 8.43 ± 0.88 | 10.25 ± 1.34 | 11.22 ± 1.08 |
| Ada-UCB ($\alpha = 1.0$) | 11.34 ± 5.84 | 12.81 ± 4.97 | 13.28 ± 1.97 | 8.70 ± 1.36 | 15.62 ± 3.38 | 10.74 ± 2.13 | 11.09 ± 1.72 | 9.93 ± 0.81 | 8.43 ± 0.88 | 10.25 ± 1.34 | 11.22 ± 1.08 |
| TiaraS | 16.30 ± 8.05 | 22.79 ± 7.73 | 12.67 ± 4.23 | 15.13 ± 2.88 | 25.40 ± 1.11 | 22.21 ± 0.00 | 12.21 ± 1.65 | 9.73 ± 1.69 | 12.60 ± 2.56 | 14.81 ± 1.08 | 16.39 ± 1.21 |
| Tiara | 26.14 ± 0.00 | 27.73 ± 0.00 | 15.92 ± 0.00 | 16.91 ± 0.00 | 25.77 ± 0.00 | 22.21 ± 0.00 | 14.06 ± 0.16 | 22.13 ± 0.00 | 16.77 ± 0.30 | 15.69 ± 0.00 | 20.33 ± 0.03 |

- **Ada- ϵ -Greedy** is a variant of ϵ -Greedy. This algorithm inserts new tags to the candidate pool $\mathcal{T}_{\text{known}}$ when new tags are found and chooses a tag from $\mathcal{T}_{\text{known}}$.
- **AdaUCB** is a variant of UCB. This algorithm inserts new tags to the candidate pool $\mathcal{T}_{\text{known}}$ when new tags are found and chooses a tag from $\mathcal{T}_{\text{known}}$.
- **TIARA-S** is a variant of TIARA that uses only the query tag for training.

4.1.2 Hyperparameters. We use $\lambda = 1$ and $\alpha = 0.01$ for TIARA and TIARA-S across all settings without further tuning. We will show that TIARA is insensitive to the choice of these hyperparameters over orders of magnitude in Section 4.5. We use 300-dimensional GloVe² trained using six billion Wikipedia 2014 + Gigaword 5 tokens for the word embeddings.

We report the performance with various hyperparameters for the baseline methods. Note that hyperparameter tuning is prohibitive in practice because of the tight API limit. If we apply hyperparameter

²<https://nlp.stanford.edu/projects/glove/>

tuning, we should use these query budgets for the main task instead. Therefore, this setting is slightly advantageous for the baseline methods.

4.2 Open Images dataset

We use the Open Images Dataset V6 [35] to construct the environment for the first testbed. Each image in this dataset has multi-class annotations, such as “cat,” “Aegean cat,” and “Pumpkin pie.” An image has 8.8 classes on average. We use these classes as tags. We utilize ResNet18 [27] trained with ImageNet for the black-box functions. Specifically, for each class c of ImageNet, we use the pre-softmax logit for c as the back-box function. Among 1000 classes of ImageNet, we use 10 classes $c = 1, 101, \dots, 901$, i.e., Tench, Black Swan, Tibetan Terrier, Tiger Beetle, Academic Gown, Cliff Dwelling, Hook, Paper Towel, Slot Machine, and Water Tower. Note that the retrieval algorithms do not know these class names, i.e., the function is a black-box. We use these class names for only evaluations and interpretations of the results. Note also that the set of tags (i.e.,



Figure 2: Safebooru Environment.³ (Left) Learning curves averaged over 10 independent runs. (Middle) Source images. (Right) Best Images found. In the first row, we show the second best image for TIARA because TIARA found the exact source image.

classes of the Open Image Dataset) differs from the set of classes of ImageNet.

We subsample 10,000 test images from the Open Images Dataset and construct an environment. When we query tag (class) t , this environment returns a random image with class t and the set of classes this image belongs to. We choose 100 random tags as the initial known tags \mathcal{T}_{ini} and set the budget to $B = 500$.

We run ten trials with different seeds. Table 2 reports the means and standard deviations of the best f for each method within B queries. The last column reports the average of ten classes. These results show that TIARA performs the best under all settings, and TIARA-S performs second best on average.

The middle panels in Figures 1 show the word clouds⁴ generated by TIARA. The size and color of a tag represent the score of the tag at the final iteration. These scores provide interpretations for the black-box functions. For example, in the case of tench, fish-related tags, such as Trout, Carp, and Milkfish, have high scores. In the black swan case, bird-related tags, such as Swan, Waterfowl, and Duck, have high scores. We stress again that TIARA does not use the ground truth class name but instead treats f as a black-box function. Even when the ground truth class name is unavailable to us, the word cloud generated by TIARA and the retrieved images indicate that f is fish-related in the first example and bird-related in the second example with significant interpretability.

4.3 Safebooru Environment

The aim of this section is to confirm that TIARA is effective even in a completely different domain. We use Safebooru as a testbed.

We use a dump retrieved on June 7, 2019⁵. Each image has 34 tags on average, such as “smile,” “long hair,” and “blonde hair.” We use illustration2vec [52] to construct the black-box functions.

We use the latest 100,000 images from this dataset and construct an environment. As a result of the broken links, this environment contains 81,517 images in total. When we query tag t , this environment returns a random image with tag t and the set of tags this image belongs to. We choose 100 random tags as the initial known tags \mathcal{T}_{ini} and set the budget to $B = 500$.

We conduct semantically similar image search experiments. We emphasize that although this environment does not provide an official image-based search, our algorithm enables such an image-based search. Given a source image s , we compute 4096-dimensional embedding $\mathbf{v}_s \in \mathbb{R}^{4096}$ using the pre-trained illustration2vec model⁶. We use the Gaussian kernel between the embeddings of an input image x and the source image as the black-box function, i.e., $f(x) = \exp(-\|\mathbf{v}_s - \mathbf{v}_x\|^2 / \sigma^2)$, where $\mathbf{v}_x \in \mathbb{R}^{4096}$ is the embedding of input x computed using the illustration2vec model, and σ is the bandwidth of the Gaussian kernel. We set $\sigma = 100$ during this experiment. The more semantically similar the input image is to the source, the higher the value taken by this function.

We use two illustrations shown in Figure 2 (middle) as the source images. The right panels in Figures 2 show images retrieved by the methods. The first source image is in the image database, and TIARA succeeds in retrieving the same image from the database. We show the second best image retrieved by TIARA in this case. Even the second best image is semantically similar to the source image, i.e., it depicts the same character with cherry blossom motifs, and has a higher objective value than the images retrieved by the baseline methods. For the second case, the source image is not in

³The source images are used under the attribution, non-commercial, and no-derivative piapro license. The other images are used with artist permission.

⁴https://github.com/amueller/word_cloud

⁵<https://www.kaggle.com/alamson/safebooru>

⁶<https://github.com/rezoo/illustration2vec>



Figure 3: Flickr Environment. (Top) Black swan. (Bottom) Similar image retrieval. (Left) Learning curves averaged over ten independent runs. (Top middle) Word cloud generated by TIARA. (Bottom middle) Source image. (Right) Best images found during each iteration.

the database. Although not exactly the same, TIARA succeeds in retrieving a semantically similar image with the same characters and motifs, i.e., horn, rabbit, and costume.

These results show the flexibility of our framework such that TIARA can be applied to not only photo-like image databases but also illustration-like image databases using appropriate black-box functions. e.g., the illustration2vec model.

4.4 Flickr Environment

The aim of this section is to confirm that TIARA is effective and readily applicable to real-world environments. We use the online Flickr environment in operation as a testbed. We also use ResNet18 trained with ImageNet to construct the black-box functions. We implement oracle \mathcal{O} by combining `flickr.photos.search` and `flickr.tags.getListPhoto` APIs. We use the `license='9,10'` option such that it returns only public domain or CC0 images. We also set the budget to $B = 500$. Because Flickr contains as many as *ten billion* images, it is challenging to find relevant images within $B = 500$ queries. Each image has 15 tags on average, including “summer,” “water,” and “sea.”

We conduct two experiments in this environment. First, we conduct the same experiment as in the open image environment. The top row of Figure 3 shows the result for the black swan class. Compared to Figure 1, TIARA in this environment learns more slowly than in the open image environment. We hypothesis that this is because the Flickr environment contains noisy tags annotated by users, which are occasionally described in foreign languages, whereas the open image environment contains only clean tags that indicate solid categories judged by annotators. Nevertheless, TIARA succeeds in retrieving black swan images in several hundred queries.

Second, we conduct similar image retrieval experiments as in the safebooru environment. We use the output of the penultimate layer of the pre-trained ResNet18 for the image embedding. We use the Gaussian kernel between the embeddings of an input image x and the source image s as the black-box function, i.e., $f(x) = \exp(-\|v_s - v_x\|^2 / \sigma^2)$, where $v_x \in \mathbb{R}^{512}$ is the embedding of input x computed using ResNet18, and σ is the bandwidth of the Gaussian kernel. We set $\sigma = 10$ in this experiment. Figure 3 (bottom, center) shows the source image. As the right panels show, TIARA succeeds in retrieving semantically similar images depicting a sunrise over the sea.

4.5 Hyperparameter Sensitivity

We investigate the hyperparameter sensitivity of TIARA in this section. TIARA has two hyperparameters, i.e., exploration coefficient α and regularization coefficient λ . TIARA also has a choice of word embeddings. It is crucial to ensure stability with respect to the hyperparameter choices because tuning the hyperparameters in a new environment is difficult with tight API budgets. The left panel of Figure 4 shows the average performance of TIARA in the Open Image Dataset environment with various values of α maintaining $\lambda = 1.0$ (default value). The x -axis is plotted on a log scale. We show the performances of TIARA-S, UCB, and Random for reference. In this plot, we do not tune the hyperparameters of TIARA-S accordingly to maintain the conciseness of the plot. This plot shows that TIARA is stable with respect to α over several orders of magnitude. The middle panel of Figure 4 shows the average performance with various values of λ maintaining $\alpha = 0.01$ (default value). The x -axis is plotted on a log scale. The result shows that TIARA is stable with respect to λ over several orders of magnitude. The right panel of Figure 4 shows the average performance with various dimensions

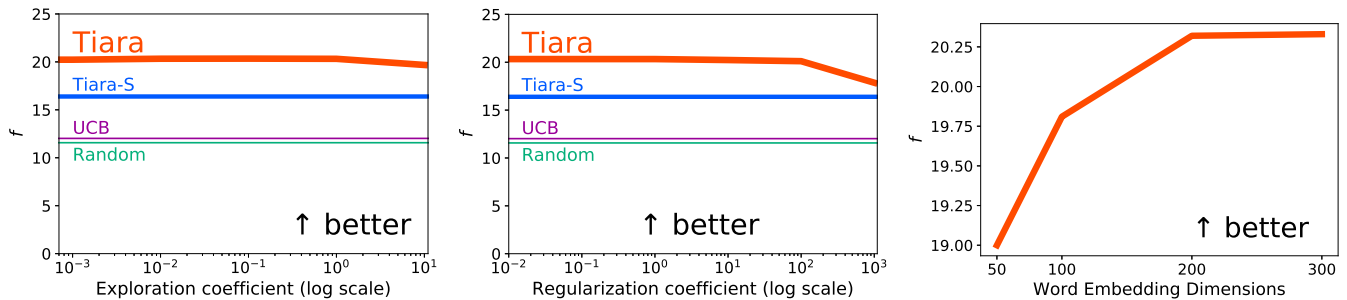


Figure 4: Hyperparameter sensitivity. As the default settings, we set the exploration coefficient as $\alpha = 0.01$, regularization coefficient as $\lambda = 1.0$, and word embedding as 300-dimensional GloVe and change one configuration in each panel. These results show that TIARA is stable with respect to the hyperparameter choices across several orders of magnitude.

of the GloVe word embedding maintaining $\alpha = 0.01$ and $\lambda = 1.0$ (default values). This result shows that TIARA performs better with higher dimensional embeddings, i.e., more expressive embeddings. It also indicates that TIARA indeed utilizes the word embedding geometry. We recommend using high-dimensional and strong word embeddings for TIARA.

5 RELATED WORK

5.1 Image Retrieval

Image retrieval has been studied for decades. The main concern is how to retrieve relevant images [4, 30, 38] with efficiency [36, 40]. Deep neural networks have been preferred for image retrieval in recent years because they can extract rich features including texture [18], style [23], and semantics [4]. Hash-based methods have also been extensively studied for their efficiency [36, 40]. In addition, many extensions have been studied, such as multi-modal image search [9, 32] and contextual retrieval [61]. We investigated the image retrieval problem from a different perspective. Specifically, we consider how an outsider can retrieve desirable images from an *external* image database with as few queries as possible.

5.2 Web Crawling

Our problem setting can be seen as a crawling problem. Developing efficient web crawlers is a long-standing problem in the literature [11, 12, 14, 15, 20, 51]. In particular, focused crawling [5, 12, 26, 29, 43] is relevant to our problem setting. Focused crawling aims to efficiently gather relevant pages by skipping irrelevant pages.

However, there are several differences between our study and existing focused crawling. First, focused crawling is used to search web pages by following WWW hyperlinks, whereas we search for *images* from an *external database* utilizing the tag search oracle. Thus, the existing crawling methods are not directly applicable to our setting. Second, we assume the API limit is extremely tight, e.g., 500, whereas existing focused crawlers typically visit hundreds of thousands of pages. Third, existing focused crawlers focus on retrieving pages of a specific topic [12, 43], popular pages [5], pages with structured data [44], or hidden web pages [7]. By contrast, deep convolutional neural networks in our setting realize rich vision applications, as we show in the experiments. These applications

are qualitatively different from the existing focused crawlers and are valuable in their own right.

5.3 Private Recommender Systems

Private recommender systems [53] aim to build a fair recommender system when the service provider does not offer a fair system. Our problem can also be seen as constructing a private recommender system when the black-box function is a recommendation score. There are several differences between our approach and private recommender systems. First, the existing methods, PRIVATE RANK and PRIVATE WALK, assume the use of an item recommendation oracle, which is unavailable under our setting. Second, TIARA considers content-based retrieval, whereas PRIVATE RANK and PRIVATE WALK mainly focus on a collaborative recommendation scenario. Third, our application is not limited to fairness, and we showed promising applications, including semantically similar image retrieval.

6 DISCUSSION AND LIMITATIONS

Although we use deep neural networks as a black-box function, our framework and the proposed method are not limited to deep neural networks. For example, we can use human judgment as the black-box function, i.e., when we evaluate $f(x)$, we actually ask a human viewer how much he/she likes image x . Because our proposed method requires several hundred evaluations for a single run, the current method is too inefficient for human-in-the-loop experiments. More efficient methods are important for developing such intriguing applications.

In this work, we have assumed that f is a well-behaved function. It may be interesting to use buggy f functions. For example, when we debug a deep neural network model f , applying TIARA to f may reveal what f has already learned and has not yet. We hypothesize that there are much room for further applications of TIARA.

We have assumed that the only way to access the image database is the tag search oracle \mathcal{O} . Although our method is general owing to this formulation, many other APIs are available in real applications, such as user-based searches, popularity ranking, and collaborative filtering recommendations. Utilizing richer information to enhance performance and query efficiency is important in practice. At the other extreme, dropping the assumptions of tag affinity and tag

search APIs to make the method applicable to broader databases is also an intriguing direction.

Although we have focused on image retrieval problems, our formulation is not limited to such problems and can be applied to other domains such as music, document, and video retrieval problems. Exploring further applications of our framework is left as future work.

7 CONCLUSION

In this paper, we formulated the problem of optimal image retrieval with respect to a given black-box function from an external image database. This problem enables each user to retrieve their preferable images from the Internet, even if the image database does not provide such features. We combined a bandit formulation with pre-trained word embeddings and proposed an effective retrieval algorithm called TIARA. Finally, we confirmed the effectiveness of TIARA using three environments, including an online Flickr environment.

ACKNOWLEDGMENTS

This work was supported by JSPS KAKENHI GrantNumber 21J22490. We are grateful to ARAM, トマリ, でんろく, Qie_, 亀村江龍, こんべ伊藤, 小本呆毛, and ミテデー for allowing the use of their wonderful illustrations.

REFERENCES

- [1] Noga Alon, Nicolò Cesa-Bianchi, Claudio Gentile, and Yishay Mansour. From bandits to experts: A tale of domination and independence. In *NeurIPS*, pages 1610–1618, 2013.
- [2] Sanjeev Arora, Yingyu Liang, and Tengyu Ma. A simple but tough-to-beat baseline for sentence embeddings. In *ICLR*, 2017.
- [3] Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Mach. Learn.*, 47(2-3):235–256, 2002.
- [4] Artem Babenko and Victor S. Lempitsky. Aggregating local deep features for image retrieval. In *ICCV*, pages 1269–1277, 2015.
- [5] Ricardo A. Baeza-Yates, Carlos Castillo, Mauricio Marín, and M. Andrea Rodríguez. Crawling a country: better strategies than breadth-first for web page ordering. In *WWW*, pages 864–872, 2005.
- [6] Ricardo Baptista and Matthias Poloczek. Bayesian optimization of combinatorial structures. In *ICML*, pages 471–480, 2018.
- [7] Luciano Barbosa and Juliana Freire. An adaptive crawler for locating hidden-web entry points. In *WWW*, pages 441–450, 2007.
- [8] Sean Bell and Kavita Bala. Learning visual similarity for product design with convolutional neural networks. *ACM Trans. Graph.*, 34(4):98:1–98:10, 2015.
- [9] Yue Cao, Mingsheng Long, Jianmin Wang, Qiang Yang, and Philip S. Yu. Deep visual-semantic hashing for cross-modal retrieval. In *KDD*, pages 1445–1454, 2016.
- [10] Alexandra Carpentier and Michal Valko. Extreme bandits. In *NeurIPS*, pages 1089–1097, 2014.
- [11] Carlos Castillo. Effective web crawling. *SIGIR Forum*, 39(1):55–56, 2005.
- [12] Soumen Chakrabarti, Martin van den Berg, and Byron Dom. Focused crawling: A new approach to topic-specific web resource discovery. *Comput. Networks*, 31(11-16):1623–1640, 1999.
- [13] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, pages 1597–1607, 2020.
- [14] Junghoo Cho and Hector Garcia-Molina. Effective page refresh policies for web crawlers. *ACM Trans. Database Syst.*, 28(4):390–426, 2003.
- [15] Junghoo Cho, Hector Garcia-Molina, and Lawrence Page. Efficient crawling through URL ordering. *Comput. Networks*, 30(1-7):161–172, 1998.
- [16] Konstantina Christakopoulou, Filip Radlinski, and Katja Hofmann. Towards conversational recommender systems. In *KDD*, pages 815–824, 2016.
- [17] Vincent A. Cicirello and Stephen F. Smith. The max K -armed bandit: A new model of exploration applied to search heuristic selection. In *AAAI*, pages 1355–1361, 2005.
- [18] Mircea Cimpoi, Subhansu Maji, and Andrea Vedaldi. Deep filter banks for texture recognition and segmentation. In *CVPR*, pages 3828–3836, 2015.
- [19] Thomas Desautels, Andreas Krause, and Joel W. Burdick. Parallelizing exploration-exploitation tradeoffs with gaussian process bandit optimization. In *ICML*, 2012.
- [20] Michelangelo Diligenti, Frans Coetzee, Steve Lawrence, C. Lee Giles, and Marco Gori. Focused crawling using context graphs. In *VLDB*, pages 527–534, 2000.
- [21] Dumitru Erhan, Yoshua Bengio, Aaron Courville, and Pascal Vincent. Visualizing higher-layer features of a deep network. *University of Montreal*, 1341(3):1, 2009.
- [22] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, pages 1126–1135, 2017.
- [23] Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *CVPR*, pages 2414–2423, 2016.
- [24] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *ICLR*, 2015.
- [25] Albert Gordo, Jon Almazán, Jérôme Revaud, and Diane Larlus. Deep image retrieval: Learning global representations for image search. In *ECCV*, pages 241–257, 2016.
- [26] Ziyu Guan, Can Wang, Chun Chen, Jiajun Bu, and Junfeng Wang. Guide focused crawler efficiently and effectively using on-line topical importance estimation. In *SIGIR*, pages 757–758, 2008.
- [27] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.
- [28] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross B. Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, pages 9726–9735, 2020.
- [29] Judy Johnson, Kostas Tsioutsoulis, and C. Lee Giles. Evolving strategies for focused web crawling. In *ICML*, pages 298–305, 2003.
- [30] Makoto P. Kato, Hiroaki Ohshima, Satoshi Oyama, and Katsumi Tanaka. Can social tagging improve web image search? In *WISE*, pages 235–249, 2008.
- [31] Tomáš Kocák, Gergely Neu, Michal Valko, and Rémi Munos. Efficient learning by implicit exploration in bandit problems with side observations. In *NeurIPS*, pages 613–621, 2014.
- [32] Saeid Balaneshin Kordan and Alexander Kotov. Deep neural architecture for multi-modal retrieval based on joint embedding space for text and images. In *WSDM*, pages 28–36, 2018.
- [33] Akshay Krishnamurthy, Alekh Agarwal, and Miroslav Dudík. Contextual semibandits via supervised learning oracles. In *NeurIPS*, pages 2388–2396, 2016.
- [34] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, pages 1106–1114, 2012.
- [35] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, Tom Duerig, and Vittorio Ferrari. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *IJCV*, 2020.
- [36] Hanjiang Lai, Yan Pan, Ye Liu, and Shuicheng Yan. Simultaneous feature learning and hash coding with deep neural networks. In *CVPR*, pages 3270–3278, 2015.
- [37] Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- [38] Luis A. Leiva, Mauricio Villegas, and Roberto Paredes. Query refinement suggestion in multimodal image retrieval with relevance feedback. In *ICMI*, pages 311–314, 2011.
- [39] Lihong Li, Wei Chu, John Langford, and Robert E. Schapire. A contextual-bandit approach to personalized news article recommendation. In *WWW*, pages 661–670, 2010.
- [40] Haomiao Liu, Ruiping Wang, Shiguang Shan, and Xilin Chen. Deep supervised hashing for fast image retrieval. In *CVPR*, pages 2064–2072, 2016.
- [41] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *ICLR*, 2018.
- [42] Shie Mannor and Ohad Shamir. From bandits to experts: On the value of side-observations. In *NeurIPS*, pages 684–692, 2011.
- [43] Andrew McCallum, Kamal Nigam, Jason Rennie, and Kristie Seymore. Automating the construction of internet portals with machine learning. *Inf. Retr.*, 3(2): 127–163, 2000.
- [44] Robert Meusel, Peter Mika, and Roi Blanco. Focused crawling for structured data. In *CIKM*, pages 1039–1048, 2014.
- [45] Tomáš Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. In *ICLR*, 2013.
- [46] Anh Mai Nguyen, Alexey Dosovitskiy, Jason Yosinski, Thomas Brox, and Jeff Clune. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. In *NeurIPS*, pages 3387–3395, 2016.
- [47] Liqiang Nie, Shuicheng Yan, Meng Wang, Richang Hong, and Tat-Seng Chua. Harvesting visual concepts for image search with complex queries. In *MM*, pages 59–68, 2012.
- [48] Wei Niu, James Caverlee, and Haokai Lu. Neural personalized ranking for image recommendation. In *WSDM*, pages 423–431, 2018.
- [49] ChangYong Oh, Jakub M. Tomczak, Efstratios Gavves, and Max Welling. Combinatorial bayesian optimization using the graph cartesian product. In *NeurIPS*, pages 2910–2920, 2019.

- [50] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. Glove: Global vectors for word representation. In *EMNLP*, pages 1532–1543, 2014.
- [51] Kien Pham, Aécio S. R. Santos, and Juliana Freire. Bootstrapping domain-specific content discovery on the web. In *WWW*, pages 1476–1486, 2019.
- [52] Masaki Saito and Yusuke Matsui. Illustration2vec: a semantic vector representation of illustrations. In *SIGGRAPH Asia Technical Briefs*, pages 5:1–5:4, 2015.
- [53] Ryoma Sato. Private recommender systems: How can users build their own fair recommender systems without log data? In *Proceedings of the 2022 SIAM International Conference on Data Mining, SDM, 2022*.
- [54] Dinghan Shen, Guoyin Wang, Wenlin Wang, Martin Renqiang Min, Qinliang Su, Yizhe Zhang, Chunyuan Li, Ricardo Henao, and Lawrence Carin. Baseline needs more love: On simple word-embedding-based models and associated pooling mechanisms. In *ACL*, pages 440–450, 2018.
- [55] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. In *ICLR workshop*, 2014.
- [56] Ashudeep Singh and Thorsten Joachims. Fairness of exposure in rankings. In *KDD*, pages 2219–2228, 2018.
- [57] Aleksandrs Slivkins. Introduction to multi-armed bandits. *Found. Trends Mach. Learn.*, 12(1-2):1–286, 2019.
- [58] Yueming Sun and Yi Zhang. Conversational recommender system. In *SIGIR*, pages 235–244, 2018.
- [59] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In Doina Precup and Yee Whye Teh, editors, *ICML*, pages 3319–3328, 2017.
- [60] William R Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3/4):285–294, 1933.
- [61] Xiaohui Xie, Jiaxin Mao, Yiqun Liu, Maarten de Rijke, Qingyao Ai, Yufei Huang, Min Zhang, and Shaoping Ma. Improving web image search with contextual information. In *CIKM*, pages 1683–1692, 2019.
- [62] Hao Yuan, Jiliang Tang, Xia Hu, and Shuiwang Ji. XGNN: towards model-level explanations of graph neural networks. In *KDD*, pages 430–438, 2020.