

# Counterfactual Attention Supervision

Seungtaek Choi, Haeju Park, Seung-won Hwang

Yonsei University, Korea

{hist0613, phj0225, seungwonh}@yonsei.ac.kr

**Abstract**—Neural attention mechanism has been used as a form of explanation for model behavior. Users can either passively consume explanation, or actively disagree with explanation then supervise attention into more proper values (attention supervision). Though attention supervision was shown to be effective in some tasks, we find the existing attention supervision is biased, for which we propose to augment counterfactual observations to debias and contribute to accuracy gains. To this end, we propose a counterfactual method to estimate such missing observations and debias the existing supervisions. We validate the effectiveness of our counterfactual supervision on widely adopted image benchmark datasets: CUFED and PEC.

## I. INTRODUCTION

Neural attention mechanism has gained interests, due to its contribution towards enhancing both accuracy and explainability. By generating a heatmap over attended regions [1] or highlighting a word of importance [2], the decision of the underlying model can be explained in a human interpretable manner. However, such work treats attention, only as a byproduct of prediction or latent variables for explanation [3], [4], while attention coefficients can also be considered as output variables, which can be human supervised.

We study the latter problem of **attention supervision (AS)**. Existing work suggests that, when such explanation coincides with human perception, accuracy also improves [5]–[10]. We illustrate our problem with an image attention supervision scenario for event type annotation [11].

Specifically, given a folder of unannotated personal images, our task is to predict its event type out of  $E$  types. Given the first row of images in Figure 1, the model is tasked to predict its event type **THEMENAME** of the given album. For this prediction, neural attention [4] may identify that the images of a ferris wheel and an animal highly contribute to the machine prediction. In **AS** problem, human can supervise attention, by giving a scalar importance score for each image in contexts of **THEMENAME** type. CUFED [12] is a dataset annotating such human supervisions, where human annotators are asked to give a scalar importance score for each image for the given event type: For the first row, the image of ferris wheel and zebra were annotated to be important for detecting **THEMENAME** event, with a high scalar score 1.5, shown as a bar and a number in Figure 1. Such score is low for the image of sky.

Our key claim is that: CUFED **attention supervision**  $S$  of image  $I$  is a **biased observation** toward the given event type  $y$ . This observation can be debiased if we can observe (or estimate) its counterfactuals: The **unobserved supervision** for image in event  $\tilde{y} \neq y$ . A closely related problem is

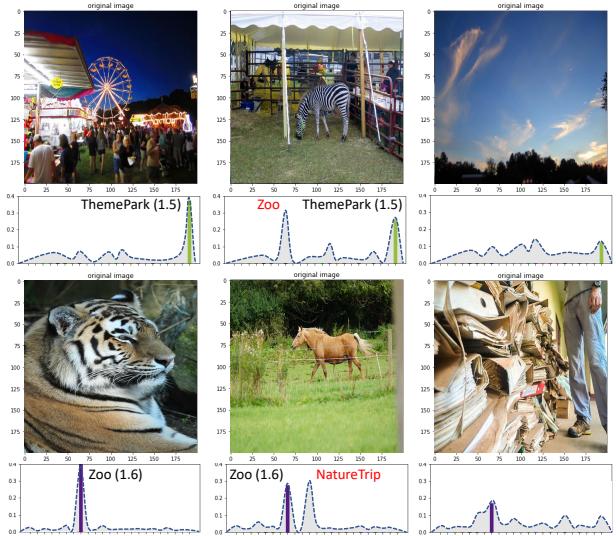


Fig. 1: Ground-truth importance score in the CUFED dataset [12] (shown as a bar) and the estimated counterfactual supervisions (shown as a distribution). The first row is from **THEMENAME** album and the second is from **Zoo** album.

obtaining an unbiased relevance estimation [13], from biased click observations to the ranking provided to the user.

One way to debias is to **collect** the counterfactual observations, or online A/B testing. In our problem setting, we may ask annotations of the same image for all  $E$  event types, which multiplies annotation overhead  $E$ -fold.

In contrast, we propose to **estimate** counterfactual observations to keep human annotation cost as low as **AS**, or offline A/B testing. That is, in Figure 1, we estimated the dotted distribution of considering the attention **distribution** for all types, where only a value shown in the bar is observed. This distributional attention view (which we name **DistAS**) allows us to realize a bias: Zebra image seems critical for detecting Zoo event, given a high observed value shown in the bar. However, the estimated annotation suggests that this image is relevant to many other types as well, and cannot contribute much to conclude the event type. In other word, its importance needs to be debiased into a lower value. This is similar in spirit with propensity weighting [14] that has been a standard approach to correct for item selection bias: interactions are biased to the documents presented at the annotation time.

Our key contribution is to leverage image semantics for propensity weighting: For example, in Figure 1, to estimate the importance of zebra in Zoo event, we can consider observed

annotations for a similar image (such as a horse in the second row) with high weights. If two given images are similar, we force the two images to have similar distribution of supervisions across event types.

We validate the effectiveness of our estimation, by *directly* comparing with human annotation on image importance, or *indirectly* by the accuracy of event type prediction. In both tasks, our proposed model, purposely built upon simple RNN and CNN models, outperforms more complex state-of-the-arts [4], [11], [12], leveraging counterfactual supervisions. Specifically, our proposed models outperform existing methods by up to 10.6% point on two personal image benchmark datasets: CUFED and PEC<sup>1</sup>.

## II. PROBLEM FORMULATION

We aim to solve the task of event type recognition for a set of unannotated images (album)  $A = \{I_1, I_2, \dots, I_T\}$ , where  $T$  is the number of images. Let  $X = \{x_1, x_2, \dots, x_T\}$  denotes the CNN features for each image of the album  $A$ , where  $x_i \in \mathbb{R}^d$  of feature dimensionality  $d$ . For event recognition, we are tasked to train a recognition function  $f : \mathbb{R}^{T \times d} \rightarrow \mathbb{R}^E$ , which predicts a correct event type  $y \in \{y_1, y_2, \dots, y_E\}$  among  $E$  event types.

Our goal is to improve the neural attention mechanism by learning attention  $\alpha \in \mathbb{R}^T$  to follow the gold importance  $S = \{S_1, S_2, \dots, S_T\}$  as closely as possible (we call **attention supervision**), which can be evaluated in the following two ways, by comparing **event-type prediction** and **event-specific image ranking** with the human annotations. For the second evaluation, we regard the attentions as an alternative of importance scoring function  $g : \mathbb{R}^{T \times d} \rightarrow \mathbb{R}^T$ , employed in the recognition function for weighting purpose  $f : \mathbb{R}^{T \times d} \xrightarrow{\alpha} \mathbb{R}^E$ .

For the sake of this discussion and without loss of generality, we will consider a decomposition of the recognition network into two functional components - an album feature extractor  $X \xrightarrow{\alpha} z$  and a decision network  $z \rightarrow \hat{y}$ . The former combines the image features  $X$  into an album representation  $z$  by weighting the image features with attention  $\alpha$ . In the latter, the album representation  $z$  is used to make event type prediction  $\hat{y}$ . Our intention is to keep this decision network as simple as possible to make the point that, with advanced attention supervision, simple models can beat more complex state-of-the-arts. We thus consider simple CNN- and RNN- based models below.

### A. CNN-Att

Dependent on the event type, the importance of images does vary and more important images should contribute more to the album representation, which can be modeled as neural attention [2], [4]. Specifically, attentions,  $\alpha = \{\alpha_1, \alpha_2, \dots, \alpha_T\}$ , are computed with a feed-forward network. We model the attentions as a probability distribution over all the images via a softmax layer. Then,  $z$  is defined by a weighted sum of image

features according to their attentions. We denote this model variant as CNN-Att. Specifically,

$$u_i = \tanh(W[x_i; x_{avg}] + b), \quad (1)$$

$$\alpha_i = \frac{\exp(u_i^\top e)}{\sum_{j=1}^T \exp(u_j^\top e)}, \quad (2)$$

$$z = \sum_{i=1}^T \alpha_i \cdot x_i, \quad (3)$$

where  $x_{avg}$  denotes the average of all the image features in the given album and  $[;]$  means concatenation of the features.  $W$ ,  $b$  and  $e$  are learnable parameters. Intuitively, attention will measure the relative importance of an image with regard to the whole album. The context vector  $e$ , representing a latent query asking for image importance for the given event.

### B. RNN-Att

Alternatively, some event type has a strong temporal dependence, such that input features are better represented as recurrent models, such as LSTM [15] and GRU [16]. We thus employ bidirectional GRU network into our attention architecture, named RNN-Att. Specifically, input images are first sorted in chronological order and fed into the BiGRU network, and hidden states of the recurrent network are used as input for the attention computation.

$$h_i = \text{BiGRU}(x_i), \quad (4)$$

$$u_i = \tanh(W[h_i; h_{avg}] + b), \quad (5)$$

$$\alpha_i = \frac{\exp(u_i^\top e)}{\sum_{j=1}^T \exp(u_j^\top e)}, \quad (6)$$

$$z = \sum_{i=1}^T \alpha_i \cdot h_i, \quad (7)$$

The decision network takes the album representation  $z$  and predicts log-probabilities over output classes ( $E$  event types).

## III. APPROACH

Our next task is to supervise such attentions for accurate prediction of both event type and importance, using public annotations, known as CUFED [12], of gold event label  $y$  and event-specific importance  $S$  for each album.

### A. Baseline: ScalarAS

Formally, we design a model to predict event type with minimal error (represented by objective function  $\mathcal{L}_{cls}$ ), but also event-specific importance (represented as  $\mathcal{L}_{ScalarAS}$ ).

First, for  $\mathcal{L}_{cls}$ , all models are trained with the classification objective, to minimize the categorical crossentropy loss  $\mathcal{L}_{cls}$  between the ground-truth  $y$  and predicted event type label  $\hat{y}$ .

$$\mathcal{L}_{cls} = \sum_{\mathcal{A}} -y \ln \hat{y}, \quad (8)$$

where  $\mathcal{A}$  denotes the entire albums in the training set.

<sup>1</sup>Our code is available at <https://github.com/hist0613/DistAS>.

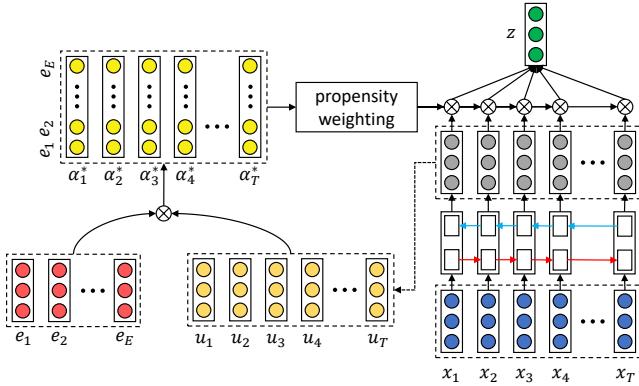


Fig. 2: The overall architecture of RNN-DistAS.

Second, for  $\mathcal{L}_{ScalarAS}$ , the objective is to ensure the distribution of attention  $\alpha$  is closer to the target distribution  $\beta$ :

$$\mathcal{L}_{ScalarAS} = \sum_{\mathcal{A}} \sum_{i=1}^T -\beta_i \log \alpha_i \quad (9)$$

Following [12], we focus on the relative importance of each image in the given album, rather than directly predicting the exact importance scores, due to the hardness of learning a reliable absolute importance. We turn the importance scores (i.e., supervisions) into a probability distribution of  $\sum_{i=1}^T \beta_i = 1$  as follows:

$$\beta_i = \frac{\exp(\lambda S_i)}{\sum_{j=1}^T \exp(\lambda S_j)}, \quad (10)$$

where  $\lambda$  is a positive hyper-parameter that controls a score contrast: When the  $\lambda$  increases, the distribution of target attention  $\beta$  becomes more skewed, guiding to attend a few of more important images.

We then set the total loss is the weighted sum of the two loss terms:  $\mathcal{L} = \mu_{cls} \cdot \mathcal{L}_{cls} + \mu_{AS} \cdot \mathcal{L}_{ScalarAS}$ , where  $\mu_{cls}$  and  $\mu_{AS}$  denote the balancing coefficients between the two terms. We apply this loss function on CNN-Att and RNN-Att respectively, and denote these variants as CNN-ScalarAS and RNN-ScalarAS.

### B. Distributional Attention Supervision (DistAS)

This section questions whether CUFED annotation  $S$  is an optimal supervision for the attention  $\alpha$ . Rather, we propose the supervision vector  $S \in \mathbb{R}^T$  should be expanded into a matrix  $S^* \in \mathbb{R}^{T \times E}$ , to annotate unobserved image importance for other event types as well. That is to say, CUFED annotation can only sparsely supervise for such matrix, by annotating  $S_{iy}^*$  for the importance for each image  $I_i$  and gold event  $y$ , namely biased towards the prediction. The same image is not considered for other types, such that  $S_{ik}^* = 0$  where  $k \neq y$ .

Now, the question is, can we augment zero entries  $S_{ik}^* = 0$  for  $k \neq y$ , with better estimates? Existing frameworks leverage the labeled data from other event types, by inventing siamese structure looking into multiple types [12], or iterative convergence [11], as implicit data augmentation. Instead, we keep structures simple and augment annotations into  $S_{ik}^*$  (replacing 0 with a counterfactual estimation), which we discuss later.

Given the expanded target supervision matrix  $S^*$ , our attention supervision goal is formally stated as follows:

$$\mathcal{L}_{DistAS} = \sum_{\mathcal{A}} \sum_{k=1}^E \sum_{i=1}^T -\beta_{ik}^* \log \alpha_{ik}^*, \quad (11)$$

where  $\beta^*$  is initialized with  $\beta_{ik}^* = \frac{\exp(\lambda S_{ik}^*)}{\sum_{j=1}^T \exp(\lambda S_{jk}^*)}$ . Note that we apply a softmax function across the images for each event type, which aims to preserve the observed ranking information within the event type. Because  $S_{ik}^*$  is zero-initialized, the softmax yields uniform distribution of  $\beta_{ik}^* = \frac{1}{T}$  for  $k \neq y$ .

To accept the expanded supervisions  $S^*$ , our attention architecture needs to be expanded to have multiple context vectors  $e_k$  as many as  $E$ , intuitively querying ‘‘important images for  $k$ -th event’’. This modification yields event-wise attention weights  $\alpha^*$  as follows:

$$\alpha_{ik}^* = \frac{\exp(u_i^\top e_k)}{\sum_{j=1}^T \exp(u_j^\top e_k)} \quad (12)$$

The overall architecture of our proposed model is presented in Figure 2.

### C. Counterfactual Supervision Estimation

From the observed importance  $S_{iy}^*$ , our goal is to estimate the unobserved importance  $S_{ik}^*$  for  $k \neq y$  at training time. The zero entries,  $S_{ik}^* = 0$ , may mean either the image is absolutely unimportant in the given event or important yet unobserved. In contrast to ScalarAS built on only the former assumption, we take the latter assumption by taking the missing supervisions  $S_{ik}^*$  as optimization variables, which can be estimated by the observed importance in other events.

Inspired by propensity weighting [14], we propose a (propensity-)weighted aggregation of observed importance for debiasing, based on the following intuition: if the two images in different events have similar image features (or, propensity), they have similar importance distributions across multiple events  $S_i^* \in \mathbb{R}^E$  (a row vector of matrix  $S^*$ ). In other words, human annotations on the given image for an unobserved event type  $S_{iy}^*$  is close to their annotation on other similar image presented for  $S_{ij}^*$ .

Formally, we set our goal as to minimize the difference between two different image similarity metrics obtained from image features and importance distributions as follows:  $\text{sim}(x_i, x_j) - \text{sim}(\beta_i^*, \beta_j^*)$ . In order to efficiently introduce such objective into existing training process, we additionally sample an album  $\tilde{A}$  whose gold event is  $\tilde{y} (\neq y)$  and build two matrices of image similarities  $M^{\text{feat}} \in \mathbb{R}^{T \times T}$  and  $M^{\text{imp}} \in \mathbb{R}^{T \times T}$  by comparing the two albums  $A$  and  $\tilde{A}$ . The  $(i, j)$ -th entry is calculated as  $M_{ij}^{\text{feat}} = \text{sim}(x_i, x_j)$  and  $M_{ij}^{\text{imp}} = \text{sim}(\beta_i^*, \beta_j^*)$ , where  $j$  denotes the index of an image in album  $\tilde{A}$ . In this work, we use cosine similarity as similarity measure, i.e.,  $\text{sim}(a, b) = \cos(a, b)$ .

Meanwhile, the above estimation provides small, yet non-zero scores for dissimilar pairs such as  $(x_{\text{elephant}}, x_{\text{mountain}})$ , generating noisy supervision. We thus redefine  $M^{\text{feat}}$  with an

introduction of threshold  $\delta$ , where an entry with value smaller than  $\delta$  becomes 0:

$$M_{ij}^{\text{feat}} = \begin{cases} \text{sim}(x_i, x_j), & \text{if } \text{sim}(x_i, x_j) > \delta, \\ 0, & \text{otherwise,} \end{cases} \quad (13)$$

where we empirically set the threshold  $\delta$  to 0.8. Such thresholding allows our model to deal with a poor estimation of feature-based image similarity.

From the two similarity matrices, we define new estimation loss  $\mathcal{L}_{\text{sim}}$  as the Frobenius norm of the error matrix  $M^{\text{feat}} - M^{\text{imp}}$ :

$$\mathcal{L}_{\text{sim}} = \|M^{\text{feat}} - M^{\text{imp}}\|_F, \quad (14)$$

where

$$\|M\|_F = \sqrt{\sum_{i \in [1, T]} \sum_{j \in [1, T]} |M_{ij}|^2} \quad (15)$$

In summary, our attention module will be jointly trained with the following two objectives:

- $\mathcal{L}_{\text{sim}}$ , estimating  $S^*$ , which is initially a sparse matrix with many unobserved importances  $S_{ij}^*$ .
- $\mathcal{L}_{\text{DistAS}}$ , supervising  $\alpha_{ik}^*$ , attached to CNN or RNN models, to follow the estimated supervision  $S_{ik}^*$ .

The entire model is trained with a new loss function:  $\mathcal{L} = \mu_{\text{cls}} \cdot \mathcal{L}_{\text{cls}} + \mu_{\text{AS}} \cdot \mathcal{L}_{\text{DistAS}} + \mu_{\text{sim}} \cdot \mathcal{L}_{\text{sim}}$ . We introduce an additional coefficient  $\mu_{\text{sim}}$  for  $\mathcal{L}_{\text{sim}}$  for balancing  $\beta^*$  estimation in the loss function. Training the expanded attention view  $\alpha^*$  with the counterfactual supervisions  $S^*$ , we name these models as CNN-DistAS and RNN-DistAS in the later experiments.

#### D. Debiased Ranking from Attention Distribution

In this section, we discuss about how we generate the debiased ranking from our attention distribution. A naive prediction of treating the maximum of the attention distribution as importance score could be inherently biased towards the observed event. We argue that better debiased ranking could be achieved by learning to discount the images, which are important in many event types, but not showing the discriminative parts, like the *zebra* image in Figure 1.

For evaluation of debiased ranking, we employ the concept of **Inverse Propensity Scoring (IPS)** [17] by giving penalty to the high frequent images across multiple event types (multiple documents). Specifically, event-specific importance, namely **relevance**  $R_i$  of the given image  $I_i$  in the album  $A$ , should be the probability of the image being relevant in *gold* event  $y$ , normalized by it being relevant in *other* events  $\tilde{y}$ , which we define as the *propensity* of the image. It can be estimated with  $\alpha^*$  as follows:

$$R_i = \frac{P(R = 1|I_i, y)}{P(R = 1|I_i, \tilde{y})} \approx \frac{S_{iy}^*}{\sum_j \text{sim}(x_i, x_j) \cdot S_{j\tilde{y}}^*} \approx \frac{\max \alpha_i^*}{1 - \max \alpha_i^*} \quad (16)$$

In this work, we treat the similarity between the two images  $(x_i, x_j)$  as a propensity score of  $x_i$  over different events. This architecture design is targeted to inference time, when we are not aware of what the gold event type is. It is not yet

guaranteed the maximum attention is of gold event  $y$ . However, by maximizing the attention score of gold event in training time, where the target supervision  $S_{iy}^*$  is initialized only at gold event, such metric could achieve correct guidance.

Finally, we obtain the album representation  $z$  according to the normalized coefficients  $r_i$  via a softmax layer:

$$r_i = \frac{\exp(R_i)}{\sum_{j=1}^T \exp(R_j)}, \quad (17)$$

$$z = \sum_{i=1}^T r_i \cdot h_i \quad (18)$$

For event-specific ranking, we use the debiased relevance score  $r_i$  as the sorting criteria for the image  $I_i$ .

## IV. EXPERIMENTS

### A. Dataset

To evaluate the effectiveness of DistAS, we conduct experiments on two public benchmark datasets: CUrator of Flickr Events Dataset (CUFED) [12] and Personal Events Collection (PEC) [18], for event recognition and event-specific ranking. Due to no available ranking annotations in PEC, we only report the result for event recognition to show the effectiveness of our counterfactual approach and debiased ranking. PEC dataset could be regarded as an extreme scenario of no human annotation.

### B. Baselines

We compare the proposed approach **DistAS** with the current state-of-the-art baselines.

- **Siamese-CNN** [12] is trained to predict the difference of importance scores between a pair of images with the piece-wise ranking loss. When evaluation, the output of CNN is used as sorting criteria.
- **Iterative-CNN-LSTM** [11] consists of three different modules: 1) CNN for image-level event recognition, 2) LSTM for album-level event recognition, and 3) Siamese networks for importance prediction from [12]. The same ResNet architecture is used as the base network in module 1 and 3. The prediction is iteratively improved by updating the output of module 1 and 2 with the importance predicted by module 3.

### C. Model Configuration

Due to the page limitation, we report the hyper-parameter settings in CUFED dataset only. The details in PEC dataset is available with our experiment codes. We use ResNet50 [19] features as the image feature  $x_i$  of dimension size 2048. The size of context vector ( $e$  and  $e_k$ ) is set to 128. For recurrent models, the size of hidden states is fixed to 512, yielding 1024 in bidirectional model. The decision network, i.e., the last fully connected layers, contains two feed-forward layers of 300-dimension with 0.2 dropout rate.

Regarding the other hyper-parameters:  $\lambda$  is empirically set to 3.0, making more clear contrast between important and unimportant images. We observe  $\mu_{\text{AS}}$  works differently in two

| Model              | Precision@K% |             |             |
|--------------------|--------------|-------------|-------------|
|                    | 10           | 20          | 30          |
| Random             | 9.0          | 19.3        | 29.8        |
| CNN-Att            | 15.1         | 28.4        | 39.9        |
| RNN-Att            | 24.4         | 39.0        | 50.3        |
| Siamese-CNN        | 28.1         | 40.4        | 49.7        |
| Iterative-CNN-LSTM | 30.0         | 41.3        | 50.7        |
| CNN-ScalarAS       | 30.9         | 48.9        | 61.0        |
| RNN-ScalarAS       | 34.4         | 50.1        | 63.5        |
| CNN-DistAS         | 36.6         | 53.0        | 63.7        |
| <b>RNN-DistAS</b>  | <b>40.6</b>  | <b>57.5</b> | <b>70.1</b> |

TABLE I: Results of event-specific ranking on CUFED.

different AS approaches: 0.2 for ScalarAS and 0.8 for DistAS. We posit that such difference stems from that the target attention  $\beta^*$  used in DistAS already contain rich information about gold event label  $y$ . For the counterfactual estimation, we observe that 0.01 for  $\mu_{sim}$  works well. There was unstable training problem when we use larger value for  $\mu_{sim}$ , such as 0.1 and 1. One possible reason is that the randomly sampled  $\tilde{A}$  introduces unnecessary training signals at the beginning of training, before learning useful ranking information.

#### D. Training Details

Following [20], the training is done following the same protocol of extracting multiple subsets from an album, where we extract 16 images (i.e.,  $T = 16$ ) over 20 times. To diminish the side-effect of such sampling, we report the average performance over 5 runs. We use Adam [21] optimizer with learning rate of 0.001. Models are trained over 50 epochs to ensure convergence of training loss with batch size of 64. All models are evaluated when showing their best ranking performance at validation set.

#### E. Direct Evaluation: Event-Specific Ranking

We begin the assessment of our model with a *direct* evaluation to show the superiority of our model. For this evaluation, we follow the protocol of [11], reporting precision@K% metric, which tells how many images of the highest predicted importance score  $\alpha_i$  are ranked in top K% images ordered by the ground truth importance.

The experimental results are shown in Table I. Our finding could be summarized as 2-fold: First, as expected, our attention supervision approaches are better able to rank images than the state-of-the-art baselines. In particular, RNN-DistAS achieves 40.6% at P@10% metric, outperforming the previous state-of-the-art Iterative-CNN-LSTM model by 10.6% point. Notably, we could observe substantial improvement even in the weakest model CNN-ScalarAS among our proposed models, achieving 7.6% at P@20% and 10.3% at P@30%, compared to Iterative-CNN-LSTM. It demonstrates the effectiveness of our problem formulation, employing the supervised attention as internal ranking function.

Second, we manifest the effectiveness of our counterfactual supervisions, particularly at P@10%. CNN-DistAS achieves the 5.7% improvement compared to CNN-ScalarAS, and



(a) Zoo event.



(b) ARCHITECTURE event.

Fig. 3: Qualitative examples of two different events. For comparison, we present the ranked list of images by 1) ground-truth importance, 2) RNN-ScalarAS, and 3) RNN-DistAS for each event. Red boxes represent the false positive images not included in the top-8 ground-truth images.

RNN-DistAS achieves 6.2% point gain over RNN-ScalarAS. It demonstrates that debiasing the importance of images, which are important at multiple events, is essential for selecting the most representative image in the given album.

For further analysis, we show qualitative examples in Figure 3. We present the top-8 ranked images by each model. As discussed, we can observe that RNN-ScalarAS incorrectly gives high scores for irrelevant images, such as the *flower* image in ARCHITECTURE (more important in NATURETRIP album). Meanwhile, our approach better highlights more discriminative image. Even when the ranked images are not optimally correlated with human-ordered images, RNN-DistAS consistently shows a reasonable ordering, such as the *tiger* image at top-1 in Zoo event, compared to *building* image of RNN-ScalarAS.

#### F. Indirect Evaluation: Album Event Recognition

The main objective of our work is to investigate the impact of counterfactual supervisions. Following the direct evaluation, here we evaluate our model in terms of their contribution to event recognition task. The results of album event recognition on the two datasets are provided in Table II. From the table, we can observe similar trends with direct evaluation, showing the strength of the debiased ranking in attention mechanism.

Our best performing model RNN-DistAS, reaching an accuracy of 75.7%, which shows a 3.4% improvement over the state-of-the-art baseline Iterative-CNN-LSTM in CUFED dataset. At the same time, RNN-DistAS achieves better performance 91.1% than Hierarchical-CNN-Att 90.1%, showing the strength of debiased ranking even in the extreme scenario of no human annotation.

| Model                    | Accuracy (%) |             |
|--------------------------|--------------|-------------|
|                          | CUFED        | PEC         |
| CNN-Att                  | 71.9         | 86.6        |
| RNN-Att                  | 72.2         | 87.1        |
| Hierarchical-CNN-Att [4] | -            | 90.1        |
| Iterative-CNN-LSTM [11]  | 72.3         | -           |
| CNN-ScalarAS             | 73.3         | -           |
| RNN-ScalarAS             | 73.7         | -           |
| CNN-DistAS               | 75.1         | 90.5        |
| RNN-DistAS               | <b>75.7</b>  | <b>91.1</b> |

TABLE II: Results of album event recognition.

## V. RELATED WORK

### A. Attention Supervision

This paper raises a bias problem in existing attention supervision, while previous literature assumes no such bias: [22] is the pioneering work that introduces the inconsistency between the attentions of human and machine in Visual Question Answering (VQA) task. Towards plausible (to human insights) attentions, [8], [10] use human attention annotations, i.e., human gaze, to supervise the attention of neural architecture in vision tasks. However, they incur expensive overheads of human annotations, such that methods for replacing human annotations are explored [7], [9], [23], [24].

Although several work propose to supervise the neural attention for each specific task, to the best of our knowledge, our work is the first to study the augmentation of counterfactual supervisions for providing improved attentions, without increasing annotation overheads on human side.

### B. Event-specific Ranking and Recognition

The goal of event recognition is to assign labels (e.g., CASUALFAMILYGATHER and BIRTHDAY) to the given image or album. With the recent advances for image understanding [19], [25], many event recognition approaches use deep learning models, such as CNN, to capture the semantic of single image (or, multiple images in the album). For example, for representing an album, to effectively combine single image features, a neural attention is introduced by [4], which we adopt as a baseline. A key distinction of our work is, we study the task of supervising such attentions, which would contribute to boosting representation quality.

## VI. CONCLUSION

In this paper, we study the problem of counterfactual attention supervision in the personal album recognition and ranking tasks. We propose to augment attention supervision by estimating the missing image importance in the counterfactual events, without additional annotation overheads. This augmented supervision can combine with simple models, improving the event-specific relevance modeling, and outperforms more sophisticated state-of-the-arts.

## ACKNOWLEDGMENTS

This work was supported by Samsung Research Funding Center of Samsung Electronics under Project Number SRFC-IT1701-01. Hwang is a corresponding author.

## REFERENCES

- [1] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio, “Show, attend and tell: Neural image caption generation with visual attention,” in *ICML*, 2015, pp. 2048–2057.
- [2] Z. Yang, D. Yang, C. Dyer, X. He, A. J. Smola, and E. H. Hovy, “Hierarchical attention networks for document classification,” in *HLT-NAAACL*, 2016, pp. 1480–1489.
- [3] L. Yu, M. Bansal, and T. Berg, “Hierarchically-attentive rnn for album summarization and storytelling,” in *EMNLP*, 2017, pp. 977–982.
- [4] C. Guo, X. Tian, and T. Mei, “Multigranular event recognition of personal photo albums,” *IEEE Transactions on Multimedia*, vol. 20, no. 7, 2018.
- [5] L. Liu, M. Utiyama, A. Finch, and E. Sumita, “Neural machine translation with supervised attention,” in *COLING*, 2016, pp. 3093–3102.
- [6] H. Mi, Z. Wang, and A. Ittycheriah, “Supervised attentions for neural machine translation,” in *EMNLP*, 2016, pp. 2283–2288.
- [7] C. Liu, J. Mao, F. Sha, and A. L. Yuille, “Attention correctness in neural image captioning,” in *AAAI*, 2017, pp. 4176–4182.
- [8] C. Gan, Y. Li, H. Li, C. Sun, and B. Gong, “Vqs: Linking segmentations to questions and answers for supervised attention in vqa and question-focused semantic segmentation,” in *Proc. IEEE Int. Conf. Comp. Vis.*, vol. 3, 2017.
- [9] T. Qiao, J. Dong, and D. Xu, “Exploring human-like attention supervision in visual question answering,” *arXiv preprint arXiv:1709.06308*, 2017.
- [10] Y. Yu, J. Choi, Y. Kim, K. Yoo, S.-H. Lee, and G. Kim, “Supervising neural attention models for video captioning by human gaze data,” in *CVPR*, 2017, pp. 2680–29.
- [11] Y. Wang, Z. Lin, X. Shen, R. Mech, G. Miller, and G. W. Cottrell, “Recognizing and curating photo albums via event-specific image importance,” in *BMVC*, 2017.
- [12] ———, “Event-specific image importance,” in *CVPR*, 2016, pp. 4810–4819.
- [13] A. Agarwal, K. Takatsu, I. Zaitsev, and T. Joachims, “A general framework for counterfactual learning-to-rank,” in *ACM Conference on Research and Development in Information Retrieval (SIGIR)*, 2019.
- [14] P. R. Rosenbaum and D. B. Rubin, “The central role of the propensity score in observational studies for causal effects,” *Biometrika*, vol. 70, no. 1, pp. 41–55, 1983.
- [15] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [16] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” *arXiv*, 2014.
- [17] R. Jagerman, H. Oosterhuis, and M. de Rijke, “To model or to intervene: A comparison of counterfactual and online learning to rank from user interactions,” 2019.
- [18] L. Bossard, M. Guillaumin, and L. Van Gool, “Event recognition in photo collections with a stopwatch hmm,” in *ICCV*, 2013.
- [19] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *CVPR*, 2016, pp. 770–778.
- [20] Z. Wu, Y. Huang, and L. Wang, “Learning representative deep features for image set analysis,” *IEEE Transactions on Multimedia*, vol. 17, no. 11, pp. 1960–1968, 2015.
- [21] D. P. Kingma and J. L. Ba, “Adam: A method for stochastic optimization,” in *Proc. 3rd Int. Conf. Learn. Representations*, 2014.
- [22] A. Das, H. Agrawal, L. Zitnick, D. Parikh, and D. Batra, “Human attention in visual question answering: Do humans and deep networks look at the same regions?” *Computer Vision and Image Understanding*, vol. 163, pp. 90–100, 2017.
- [23] Y. Zhang, J. C. Niebles, and A. Soto, “Interpretable visual question answering by visual grounding from attention supervision mining,” *arXiv preprint arXiv:1808.00265*, 2018.
- [24] Z. Wang, X. Liu, L. Chen, L. Wang, Y. Qiao, X. Xie, and C. Fowlkes, “Structured triplet learning with pos-tag guided attention for visual question answering,” *arXiv preprint arXiv:1801.07853*, 2018.
- [25] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.