

How to Become Instagram Famous: Post Popularity Prediction with Dual-Attention

Zhongping Zhang
Electrical and Computer Engineering
University of Rochester
Rochester, NY 14627
zhongping.vista@gmail.com

Tianlang Chen
Department of Computer Science
University of Rochester
Rochester, NY 14627
tchen45@cs.rochester.edu

Zheng Zhou
Department of Electrical Engineering
University at Buffalo
Buffalo, NY 14260
zzhou32@buffalo.edu

Jiaxin Li
Department of Electrical Engineering
Harbin Institute of Technology
Harbin, China, 150000
lijiaxinpp93@yeah.net

Jiebo Luo
Department of Computer Science
University of Rochester
Rochester, NY 14627
jluo@cs.rochester.edu

Abstract—With a growing number of social apps, people have become increasingly willing to share their everyday photos and events on social media platforms, such as Facebook, Instagram, and WeChat. In social media data mining, post popularity prediction has received much attention from both data scientists and psychologists. Existing research focuses more on exploring the post popularity on a population of users and including comprehensive factors such as temporal information, user connections, number of comments, and so on. However, these frameworks are not suitable for guiding a specific user to make a popular post because the attributes of this user are fixed. Therefore, previous frameworks can only answer the question “whether a post is popular” rather than “how to become famous by popular posts”. In this paper, we aim at predicting the popularity of a post for a specific user and mining the patterns behind the popularity. To this end, we first collect data from Instagram. We then design a method to figure out the user environment, representing the content that a specific user is very likely to post. Based on the relevant data, we devise a novel dual-attention model to incorporate image, caption, and user environment. The dual-attention model basically consists of two parts, explicit attention for image-caption pairs and implicit attention for user environment. A hierarchical structure is devised to concatenate the explicit attention part and implicit attention part. We conduct a series of experiments to validate the effectiveness of our model and investigate the factors that can influence the popularity. The classification results show that our model outperforms the baselines, and a statistical analysis identifies what kind of pictures or captions can help the user achieve a relatively high “likes” number.

Keywords-Popularity Prediction, Dual-Attention Model, Instagram, Social Media Data Mining

I. INTRODUCTION

Social media platforms provide their users with a good opportunity to share daily lives, emotions, and so on. Driven by data, post popularity prediction has been focused and studied widely in recent years. Researchers are able to build powerful models to predict post popularity from various aspects, such as image [1][2], textual content [3][4], time

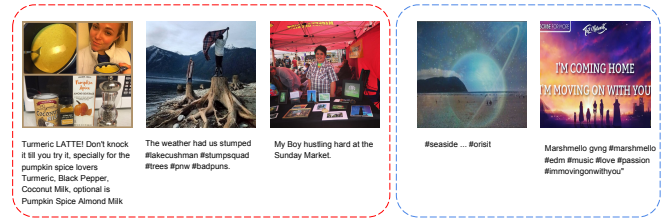


Figure 1: Example posts in Instagram. Posts in the red dotted box are popular. Posts in the blue dotted box are unpopular.

series [5], sentiment [6] or even brand information [7]. These frameworks always measure the popularity of a post from the view of the whole social media platform. For example, the authors in [1] and [3] respectively use the number of views and the forwarding number as measurements. Though these indexes can measure popularity from a big picture level, they ignore the diversity of users. For instance, a post with 100 ~ 200 views is popular for a new user. However, a post with 10000 ~ 20000 views still might not be popular for a famous star on the same platform. Under this circumstance, the above measurements are not able to reflect the post popularity for a particular user. A big shot always obtains high popularity scores while a green hand tends to be assigned a low popularity score. In practice, post popularity prediction for a particular user is significant for both companies and their customers. Companies are able to maximize their influence and provide more compelling content for their users. Customers who want to become more attractive can evaluate their posts before they upload them. Motivated by these benefits, we raise a new target in this paper: post popularity prediction for a specific user. To solve the problem, we develop a system which takes the user’s images and captions as input and then generates the popularity prediction result. To make our discussion more

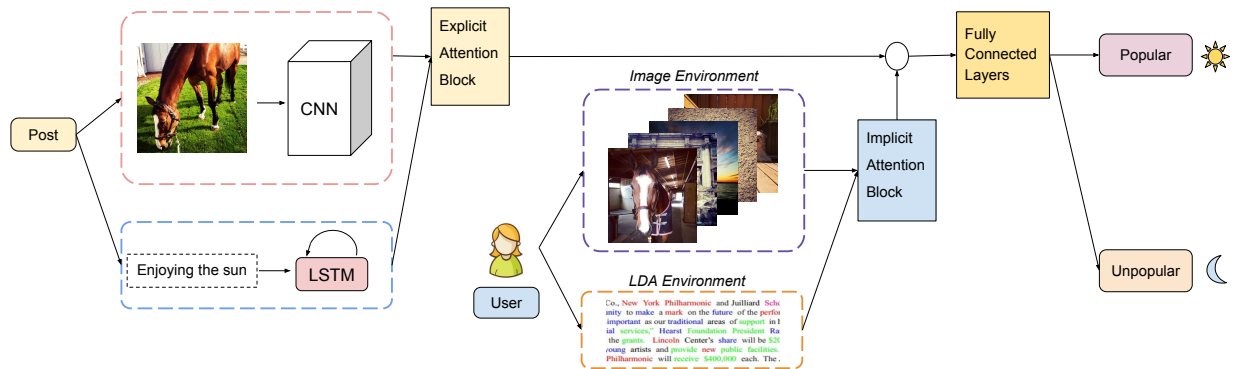


Figure 2: Overview of the dual-attention model.

straightforward, we illustrate several examples in Figure 1.

We formulate our task as a binary classification problem to classify whether a post is popular for a particular user. In this paper, a novel dual-attention model is proposed to predict the result. Concretely, the dual-attention model includes two parts: explicit attention model and implicit attention model. These two models take different levels of information as input, and then they are concatenated by a hierarchical structure. Specifically speaking, the explicit attention model is designed to generate attention weights for captions. We modify a co-attention model [8] to make it more suitable for Instagram image-caption pairs. Implicit attention model is applied to incorporate user environment. The user environment includes image environment and topic environment. Since environment does not have any explicit meaning for different positions, an implicit attention model without explicit attention mechanism is deployed here. We use a hierarchical structure to connect the explicit attention model and the implicit attention model. The structure is designed because the user environment contains higher level information than image or caption alone. Figure 2 demonstrates the framework of our model.

In this study, we collect our data from Instagram, a platform where people can share their pictures and emotions. The dataset contains 441 users and 60,785 image-caption pairs. Based on image-caption pairs, we extract the user environment and feed it into the proposed implicit attention model. Our target is not limited to predict the popularity of a post, but also explore the correlations between image, caption, and popularity. A series of experiments is therefore performed to evaluate our model and reveal the correlations.

The main contributions of our work are:

- We introduce a framework to address the problem of post popularity prediction for a specific user.
- We propose a method to calculate user environment. Compared with image and caption, the user environment is a higher level information. It can provide the model with a more comprehensive understanding of users, thus

can further improve the performance of the model.

- We develop a novel dual-attention model to predict whether a post is popular. The dual-attention model consists of two parts, an explicit attention part for image-caption pairs and an implicit attention part for user environment.
- We perform two levels of experiments on the Instagram dataset. First, we present the classification results to demonstrate the effectiveness of our framework. Next, we explore the factors which can influence the user’s post popularity.

II. RELATED WORK

Our work is mainly related to user trait pattern, popularity prediction, and attention model. In this section, we will respectively discuss the related work from these three aspects.

A. User Trait Pattern

With recent advances in social media data mining, exploring user trait behind data has become a popular research topic. Image information is extensively used to provide valuable cues for identifying user attributes. In [9], the authors use images posted on various social networks to infer the user gender. More recently, Dhir et al. [10] attempt to predict age and gender from selfie-related behavior. Besides, there are several efforts to explore the inner traits of a person, such as personality [11], interest [12], etc. Topic information is another important source for user trait prediction. One of the most influential papers in topic model is Latent Dirichlet Allocation (LDA) [13]. LDA is a generative probabilistic model which is widely used to extract topics from unlabeled documents. Based on the topic model and text messages, many interesting studies such as [14][15] are presented in recent years.

In our paper, we apply LDA topic model to construct user topic environment. We will discuss more details in Section III-C.

B. Attention Model

Attention model is first applied in English-French translation task [16]. Motivated by the success in language translation, many researchers apply attention models in image captioning [17][18][19]. To get a better representation of image and caption, co-attention model [20][8] are proposed. Since co-attention mechanism considers both the image and caption, it can generate affinity matrix which includes not only the spatial attention weights for image but also the text attention weights for caption. Considering the attention weights can be extracted from a certain layer of the model, this kind of attention mechanism can be defined as explicit attention mechanism. On the other hand, Kim et al. [21] propose an implicit attention without explicit attention parameters. They apply the structural similarity with residual learning to avoid the attention parameters, but still effectively learns the joint representation from vision and natural language.

C. Popularity Prediction

The technological and economic importance of popularity prediction motivate many researchers to notice this area [3][4][22][23][5][24][6][7][25][26]. Image is the main research direction of popularity prediction. The authors in [22][23][27][1][2] explore image popularity based on information extracted from image, like objects, image metadata and so on. Compared with image, text is another hot research area. In [3][4], the authors predict on-line message popularity by analyzing textual information. Besides image and text, some novel information such as time [5][24], sentiment [6] or even brand [7][25] are also considered to guide the prediction.

As we discussed above, previous work on popularity prediction always focuses on the big picture level and ignores the diversity of users. Furthermore, many of these frameworks attempt to incorporate extra information like friend links, user contacts, sentiments, and user tags to improve the accuracy of prediction. However, the above information is not always available. In this paper, we attempt to develop a system which can predict post popularity for a particular user based only on image-caption pairs.

III. DUAL-ATTENTION MODEL

In this section, we introduce a dual-attention model by five steps. Firstly, we start by listing some important notations to avoid ambiguity. Secondly, we present the explicit attention model in Section III-B. Then, we introduce the user environment and describe the details on how to calculate it in Section III-C. After that, the implicit attention model is proposed in Section III-D. Finally, the overall structure of the dual-attention model is described.

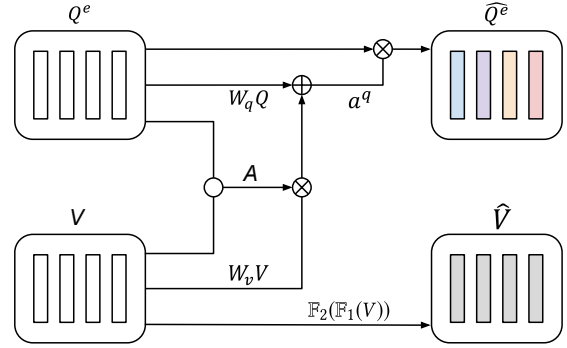


Figure 3: Illustration of the explicit attention model.

A. Notations

To ease understanding the following parts, here we list several important notations:

- $Q = \{q_1, \dots, q_T\}$ denotes a caption with T words, where q_t corresponds the onehot vector of the t -th word
- $W_{(\cdot)}$ denotes weights of different layers, we omit bias to avoid redundancy
- $\sigma_{(\cdot)}$ denotes activation functions of different layers
- $\mathbb{F}(\cdot)$ denotes fully connected layer
- \hat{V} and \hat{Q}^e denote the attended features of image and caption
- I_e and T_e denote image environment and topic environment respectively
- $F(I_e, T_e)$ and $H(q, v)$ denote joint residual function and optimal mapping respectively, these two notations are consistent with the definition in [21]

B. Explicit Attention Model

Given an image I and its corresponding caption $Q = \{q_1, \dots, q_T\}$, we first encode them into a feature vector space. ResNet-50 [28], which is pre-trained on ImageNet dataset [29], is applied here as the image encoder. We extract the image features V from the last pooling layer, whose dimension is 2048. To get word features $Q^e = \{q_1^e, \dots, q_T^e\}$, we embed the captions with a word embedding layer followed by a one-layer LSTM.

Co-attention mechanism [30][8] are commonly used in order to get a better representation of images and words. However, unlike the traditional visual question-answering (VQA) or image captioning tasks, Instagram images usually do not contain complex spatial information because most pictures uploaded by users are selfies, landscapes, posters and so on. Most of these pictures do not express complicated logical relationship nor different importance among objects. Therefore, we modify the co-attention model based on [8] to make it more suitable for our popularity prediction task.

Concretely, the attention model starts with calculating the affinity matrix A between image $V \in R^{d_1}$ and caption $Q^e \in$

$R^{d_2 \times T}$ representations

$$A = \tanh((Q^e)^T W_a V) \quad (1)$$

where $W_a \in R^{d_2 \times d_1}$ is the learning matrix. Elements in $A \in R^{T \times 1}$ are affinity scores between image and each word. According to affinity matrix A , we can further calculate the attention weights a^q via the following equations

$$\begin{aligned} H^q &= \tanh(W_q Q + A \cdot W_v V) \\ a^q &= \text{softmax}(W_h H^q) \end{aligned} \quad (2)$$

where $W_q \in R^{k \times d_2}$, $W_v \in R^{k \times d_1}$ and $W_h \in R^{k \times 1}$ are all learning matrices for the explicit attention model. k represents the last dimension of H^q , and here we manually set k as 128. Note that though we only obtain the attention weights for captions, image information is still involved during the process of calculating a^q .

Finally, the new representations of image and caption are

$$\begin{aligned} \hat{V} &= \mathbb{F}_2(\mathbb{F}_1(V)) \\ \hat{Q}^e &= \sum_{t=1}^T a_t^q q_t^e \end{aligned} \quad (3)$$

$\mathbb{F}(\cdot)$ means fully connected layer. We apply $\mathbb{F}(\cdot)$ to map the image features V into the same dimension as caption features \hat{Q}^e .

We demonstrate the structure of our attention model in Figure 3. Since the parameters of attention weights are explicitly propagated in this model, this model is named as explicit attention model to distinguish from the implicit attention model in Section III-D.

C. User Environment

In most cases, post popularity is not only influenced by its corresponding image and caption but also rely on the user who makes the post. For example, a person whose picture wall is full of landscapes, upload a selfie one day. The selfie is very likely to get a high number of ‘‘likes’’. On the contrary, if the user is selfie-addicted and uploads selfie every day, a new selfie is less likely to be popular because his friends have got used to it. Motivated by this circumstance, we introduce the concept of user environment to further improve our model.

We utilize user environment to indicate the content that the user is very likely to post. Therefore, we introduce the average value of user features to represent the environment. Image environment I_v is directly defined as the mean value of the deep-level image features V . With respect to topic environment T_e , Latent Dirichlet Allocation (LDA) [13] is applied to assign a topic feature to each caption. LDA is a generative statistical model of corpus. It assumes that documents have several random latent topics, and each topic can be characterized by a distribution over words. In our system, we set the number of topics as 400 which means the LDA feature of each caption is a 400 dimension vector.

Similar with image features, we also use the mean value of LDA features to represent topic environment.

D. Implicit Attention Model

Given image environment I_e and topic environment T_e , the most direct method to incorporate them is using fully connected layers followed by a concatenation layer. Considering the structure of the fully connected layer cannot highlight important positions on the environment features, we apply an attention model as an alternative choice in the environment-encoding process. Different from common image and text fusion target, it is a challenging task to explicitly express what the elements of environment features stand for. Therefore we devise an implicit attention model motivated by [21]. (Consistent with [21], we still use F to denote joint residual function, and use H to denote optimal mapping.)

Different from explicit attention models, the attention parameters of implicit attention model are hidden in the element-wise multiplication layer. We present our implicit attention model in Figure 4. As the figure shows, user environment variables I_e and T_e are fed into a fully connected layer respectively in our model. The joint residual function is given by

$$F(I_e, T_e) = \sigma(W_i I_e) \odot \sigma(W_t T_e) \quad (4)$$

where σ is *Relu* and \odot is element-wise multiplication. W_i and W_t is used for encoding I_e, T_e .

Given the joint residual function, optimal mapping $H(q, v)$ is predicted by

$$H(q, v) = W_{i2} I_e + W_{t2} T_e + F(I_e, T_e) \quad (5)$$

W_{i2} and W_{t2} are shortcut for environment and both of them are encoded into the same feature dimension with $H(q, v)$. After calculating $H(q, v)$, the final representation of environment e_w is obtained by $\mathbb{F}(H(q, v))$.

E. Overall Structure for Predicting Results

After computing the attended features \hat{V} and \hat{Q}^e , environment representation e_w , we predict the final answer through a hierarchical structure as the following equations:

$$\begin{aligned} h_w &= (\hat{V} + \hat{Q}^e) \\ O_1 &= \sigma_1(W_1[h_w, e_w]) \\ O_2 &= \sigma_2(W_2 O_1) \end{aligned} \quad (6)$$

where σ_1 is *Relu* and σ_2 is *Sigmoid* function. $[\cdot]$ indicates concatenation between two tensors. Since we treat the popularity prediction problem as a binary classification task, the loss function is:

$$\frac{1}{N} \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)] \quad (7)$$

where y_i corresponds to ground truth labels and \hat{y}_i indicates predicted labels.

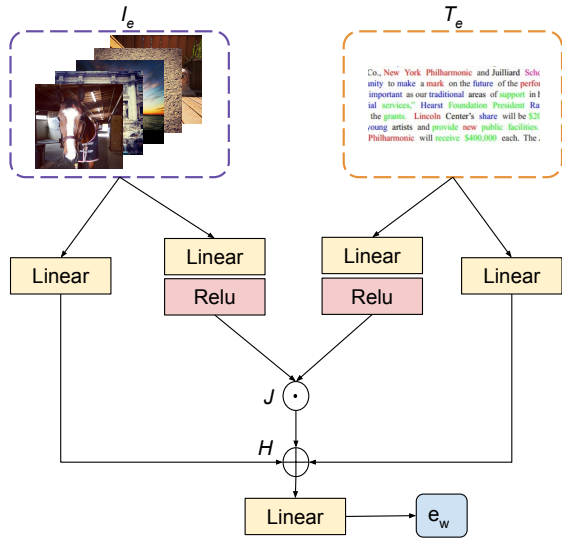


Figure 4: Structure of the implicit attention model.

IV. EXPERIMENT

A. Dataset

There is no public dataset for post popularity prediction. As a result, we collect our data by crawling on Instagram. The dataset we construct contains 441 users and their 60785 image-caption pairs, along with the corresponding number of “likes”. We choose the number of “likes” as the index to measure popularity. In order to consider the popularity for each user, we select top 25% posts (according to the “likes” number among each user’s posts) as positive samples and bottom 25% posts as negative samples. We randomly select 20% of them as the test set. Besides, we randomly split 10% of training set as the validation set to decide hyper-parameters. In the end, there are 21874 image-caption pairs for training, 2430 image-caption pairs for validation, and 6064 image-caption pairs for testing. The ratio between positive and negative samples is 1:1.

B. Classification Evaluation

In this section, we conduct classification experiments to evaluate the effectiveness of our proposed model. Considering that our system takes image-caption pairs as input, we choose the following image-caption fusion frameworks as our baselines:

- **Single Visual.** The input of Visual model only includes images. We use ResNet-50 which is pretrained on ImageNet to extract image features and feed them into the fully connected layers.
- **Single Textual** Textual features are first extracted by word embedding layer and one-layer LSTM, then fed into the fully connected layers.
- **Early Fusion.** The image and textual information are concatenated in feature level.

- **Late Fusion.** The image and textual features are fused until the last layer of the model. In another word, the final prediction score can be considered as the average value of visual prediction score and textual prediction score.
- **CCR.** CCR [31] denotes Cross-modality Consistent Regression Model. It applies KL divergence to measure the consistency between different modality features and concatenated features.
- **Similarity.** Similar with [32], we use the inner product (cosine similarity) between image and caption as their representation and feed it into the following layers.

In this paragraph, to avoid confusion, we will explain why we do not choose recent popularity prediction frameworks as our baselines. Typical popularity frameworks [1][25][33] focus on involving more useful information for prediction. Take Mazloom et al. [33] as an example, the authors introduce a three-dimensional tensor which incorporates the user category, item category, and context category. Since the input of their model is a large matrix (the three-dimensional tensor) which already contains obvious and comprehensive information of users, they apply a modified Factorization Machine (FM) to generate prediction results. Similarly, the authors in [1] and [25] use Support Vector Regression (SVR) as their prediction model. However, in this paper, the dual-attention model takes the raw image-caption pairs as input. We assume that only the image-caption pairs are available because we aim at predicting the post popularity for particular users. As is known to all, it is meaningless and unfair to compare traditional algorithms like SVR, FM with neural network methods directly on raw data. Therefore, we choose the baselines mentioned above rather than SVR, FM or other traditional algorithms.

For all the baselines and our proposed model, we apply ResNet-50 as the image encoder and one-layer LSTM as the textual encoder. The 2048-dimension image features are extracted from the last pooling layer of ResNet-50. To conduct a fair comparison, we set the dimension of word-embedding features and LSTM hidden state to 512 for all frameworks. During our training process, Adam optimization is used with a learning rate of 0.001 for the first two epochs and with a learning rate of 0.0001 for the following epochs. The size of mini-batch is set to 128.

We demonstrate the quantitative results of our experiments in Table I. The performance of different models is evaluated by four metrics: precision, recall, F-measure, and accuracy. We first compare the baselines with Explicit Attention model. As the table shows, Explicit Attention model can achieve better results under F-measure and accuracy than the other baselines. Although CCR and Late Fusion achieve relatively higher scores in recall and precision respectively, Explicit Attention model obtains a better trade-off under all the metrics. To further improve the model, we include user environment and implicit attention model to construct Dual-attention model. Since the user environment calculation does

Table I: Comparison of accuracy, precision, recall and F-score

	Precision	Recall	F-measure	Accuracy
Visual	58.61	59.76	59.18	58.34
Textual	65.09	72.48	68.59	66.46
Early Fusion	66.56	71.70	69.02	67.49
Late Fusion	67.44	66.71	67.07	66.90
CCR [31]	63.90	74.99	69.01	65.96
Similarity [32]	65.27	71.02	68.02	66.26
Explicit Attention(ours)	67.25	71.64	69.38	68.05
Dual-Attention(ours)	69.91	75.45	72.58	71.19

Table II: Comparison of accuracy, precision, recall and F-score in the ablation study

	Precision	Recall	F-measure	Accuracy
Early Fusion	66.56	71.70	69.02	67.49
E-attn	67.25	71.64	69.38	68.05
Env	67.26	69.55	68.39	67.51
Env+I-attn	69.85	67.00	68.40	68.72
E-attn+Env	70.10	72.39	71.23	70.45
E-attn+Env+I-attn	69.91	75.45	72.58	71.19

not rely on any extra information, (our user environments are extracted from users’ images and captions), we compare Dual-attention model together with the other models. The results show that Dual-attention model can take one step further based on Explicit Attention model. Almost all metrics of Dual-attention model can exceed 70%.

An interesting finding in our experiments is that: Single Textual model performs much better than Single Visual model. For F-measure or Recall, it even achieves better results than some fusion models like Similarity, Late Fusion. Based on this result, we infer that caption information plays a more reliable role than image information on post popularity prediction task.

C. Ablation Study

In order to achieve a better understanding of the effectiveness of our proposed models, we perform ablation study and present the results in Table II. In Table II, E-attn is the abbreviation for explicit attention and I-attn is the abbreviation for implicit attention. Env indicates the model that concatenates user environment directly by a hierarchical structure.

For our proposed models, since image and caption are fused in feature level, we use the Early Fusion model as our baseline. The ablation study results show that both E-attn and Env could achieve better improvement in accuracy and precision compared with Early Fusion. Besides, the performance of model improves as the number of additive structures increases. E-attn+Env model begins to surpass Early Fusion from all indexes. And finally, E-attn+Env+I-attn, namely Dual-attention model achieve best results among all combinations.

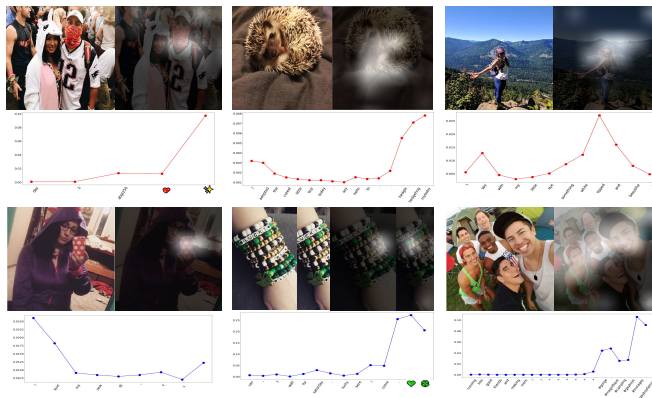


Figure 5: Visualization of the image attention maps and word attention weights. The first row corresponds to positive examples, where we demonstrate the word attention weights by the red line. The second row corresponds to negative examples, where we demonstrate the word attention weights by the blue line.

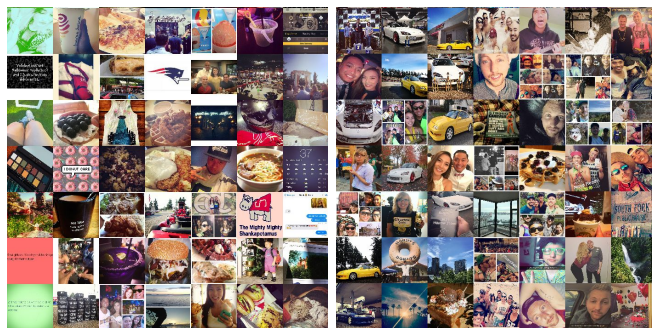


Figure 6: The clustering results for popular and unpopular images. Unpopular images are shown on the left hand side. Popular images are shown on the right hand side.

Table III: Clustering categories and their corresponding ratio R

Name	ratio R
Photos of Daily Life	0.475
Meal and Drink	0.696
Small Group Photo	0.396
Group Photo	0.393
Poster	0.760
Picture with a caption on	0.627
Landscape Photo	0.602
Text Poster	0.700
Image combined by small pictures	0.492
Phone Screen Shot	0.595
Selfie	0.461
Car and Daily Commodities	0.495

D. Visualization

To provide a deeper insight into what kind of things our model tends to focus on, we randomly select three pictures respectively from positive and negative samples, then visualize their image attention maps and word attention weights in Figure 5. For image attention maps, we extract new image features from the last Convolutional layer of ResNet-50 first. The dimension of new image features is (7,7,2048). Unlike the original image features (2048 dimension vector), the new image features contain spatial information (7*7). We input the new image features into Dual-attention model part by part and generate a probability map for each image based on the probability score of each part. Therefore, different from the traditional image attention map, our image attention map actually reflect the popularity level of each zone. As shown in Figure 5, human face, hedgehog, bracelet and mobile phone get relatively high popularity score compared with the other parts of the images. Based on this phenomenon, we conclude that concrete objects tend to get high popularity scores by our proposed model.

Under each image, we plot the word attention weights to illustrate the effectiveness of our explicit attention model. We can see that emoji, hashtag, specific object or specific action always tend to obtain high attention weights. For instance, in the first and fifth plots (from upper left to bottom right), the attention value of emoji “star”, “four-leaf clover” and “heart” are much higher than the other words in the same captions. In the last plot, all hashtags obtain high attention scores. From the second and third plots, we can observe that attention weights of “hedgehog”, “hedgie”, “spy” and “tipped” increases remarkably, indicating that the model pays more attention to these words. Generally speaking, the explicit attention mechanism is able to capture keywords in captions and correlate them well with the image information.

E. Image Analysis

1) *Intuitive Feeling*: Firstly, we would like to have an intuitive feeling about the difference between popular and unpopular images. As shown in Figure 6, popular images seem more complex and always contain objects like people, selfies, and so forth. On the other hand, unpopular images are simpler. Many of them are posters, advertisements, screenshots, or foods.

2) *K-means Clustering*: To get a more solid conclusion, we implement K-means Clustering on the high-level image features V into 12 categories. For each category, the ratio of unpopular images to total images R is set as an evaluation index. In order to obtain high-quality clustering results, we apply the following training strategy:

1. Cluster the remaining images into K classes and calculate the ratio R of each class;
2. Pick out the classes whose score $|R - 0.5|$ is larger than threshold t ;

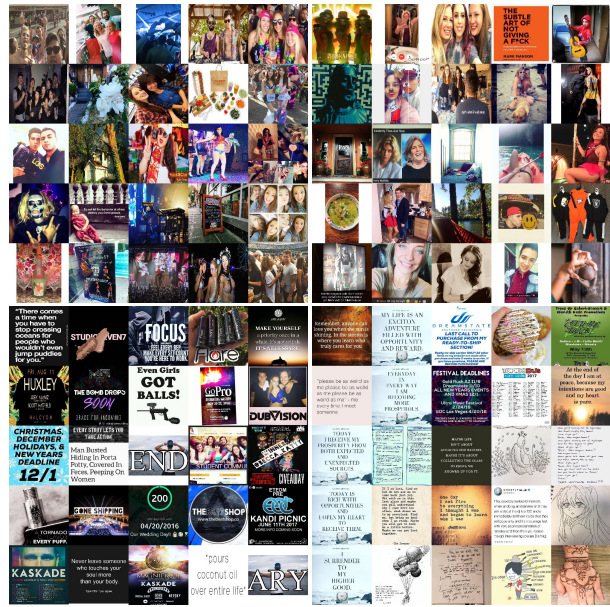


Figure 7: Selected categories by K-means clustering. From upper left to bottom right: Group Photo, Small Group Photo, Poster, Text Poster.

3. Repeat step 1 until there is no class satisfy the qualification in step 2;
4. Collect the picked out and the remaining classes as the final results.

By this strategy, the difference between popular and unpopular images of each category can be maximized. Table III shows the clustering classes and their corresponding ratio R . We can see that categories such as “poster”, “meal and drink” obtain high scores of R , which means images belong to these categories are likely to get a small number of “like”. In contrast, classes like “group photo”, “small group photo”, and “selfie” get low scores of R , which means images in these categories tend to receive a large number of “like”. This phenomenon is consistent with our intuitive conclusion in Section IV-E1. To better understand our results, we select four typical categories according to their value of R and demonstrate their pictures in Figure 7.

F. Text Statistic Analysis

Thanks to visualization, we already have an intuitive feeling about the connection between text and popularity. In this section, we analyze their correlation from a statistical perspective. We manually divide the captions into two categories: words and emojis. To avoid redundancy, we use “text” to represent the set of words and emojis. And for each category, we design the following experiment to reveal their correlations with popularity.

In this experiment, we count the frequency of occurrence for each word. Take the word “love” as an example, we go

through all the captions in positive samples and count the total occurrence number m_p of “love”. Then we perform a similar process to count the total occurrence number m_n in negative samples. Since the ratio between positive and negative samples is 1:1, we do not need to normalize m_p or m_n . In the end, the texts are ranked directly by the value of $m_p - m_n$. Before moving to the statistic analysis, we first filter out unrelated texts to avoid bias. These texts include basic symbols (“.”, “:”, “;”, ...), pronouns (“this”, “it”, “that”, ...), prepositions (“with”, “at”, “of”, ...). There are some words that are used only in special conditions. For example, “uscevents”, “paradiso” tend to appear in popular captions only from the users of USC (University of Southern California). On the other hand, “chicago”, “illinois” always appear in popular posts from the users of UIUC (University of Illinois Urbana-Champaign). Hence we also delete this kind of words. After filtering out the unrelated texts, we demonstrate the occurrence number of typical texts in Figure 8. We manually separate the words into three stages: top 25%, median 50%, and bottom 25%. In contrast, we divide the emojis into two levels: top 50% and bottom 50%. The reason is that most emojis appear in popular posts and there is not a clear boundary between unpopular emojis and the other emojis. In Figure 8 (a), words from “love” to “day” are selected from the top 25%. Based on statistic results, words that describe time (“year”, “day”, “time”), attribute (“amazing”, “beautiful”) and correlated with holiday (“festival”, “weekend”, “selfie”) are very likely show up among the top 25%. It seems that the posts include these words have higher tendency to receive “like” from other users. Words from “best” to “books” are selected from the median 50%. We observe that most nouns distribute in this range for the reason that nouns are always used to describe objective events and they do not have much sentiment or subjective opinions involved. Words from “breakfasts” to “birthday” are selected from the bottom 25%. We find that words describe food, such as “coffee”, “dinner”, “breakfast” appear frequently in this range, which is consistent with our image analysis results. Out of our expectation, “holidays”, “birthday”, “bestfriends” are in this region too. Just a guess, “holidays” is too general, it often appears in posters where people don’t always give a “like”. And when people mention “birthday”, “bestfriends” in posts, only their close friends are likely to give them a “like”. In Figure 8 (b), the first thirteenth emojis are selected from the top 50% and the others are from the bottom 50%. Since we do not find out an obvious regular variety among these emojis, we just list their statistic results and provide readers an intuitive feeling.

V. CONCLUSIONS

In this paper, we propose a dual-attention framework to address the image-caption based popularity prediction problem. Since our prediction target is for a specific user, we introduce the user environment as a high-level input to guide

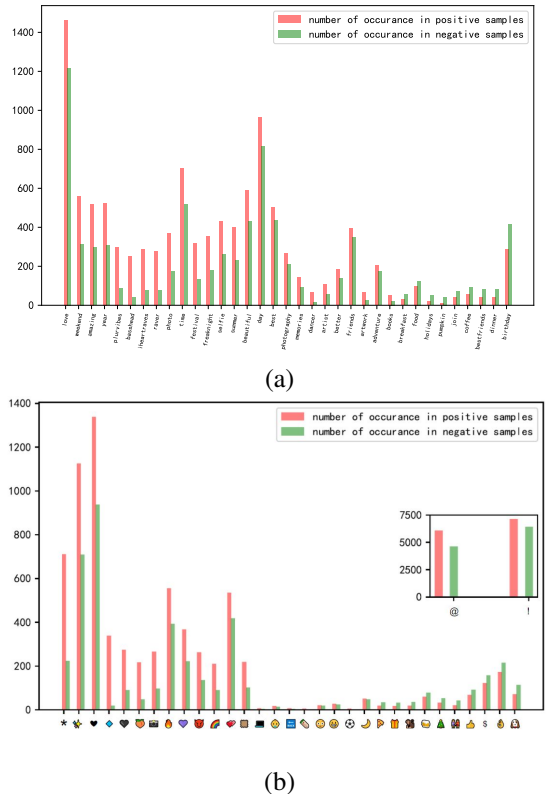


Figure 8: The occurrence number of selected texts. (a) Statistic results of words, b) Statistic results of emojis.

the classification model. User environment is incorporated by an implicit attention mechanism and image-caption pair is incorporated by an explicit attention mechanism. The classification results show that Dual-attention outperforms baselines on all measurements. Visualization results confirm that our model can clearly learn popularity words or image regions. Finally, we analysis image and textual information based on statistical results and draw conclusions about the correlation between image, caption, and popularity. In the future, we intend to develop a more efficient model to incorporate the user environment. Furthermore, a more comprehensive user profile by fusing locations, seasons, etc., can also be considered.

REFERENCES

- [1] F. Gelli, T. Uricchio, M. Bertini, A. Del Bimbo, and S.-F. Chang, “Image popularity prediction in social media using sentiment and context features,” in *Proceedings of the 23rd ACM international conference on Multimedia*. ACM, 2015, pp. 907–910.
- [2] L. C. Totti, F. A. Costa, S. Avila, E. Valle, W. Meira Jr, and V. Almeida, “The impact of visual attributes on online image diffusion,” in *Proceedings of the 2014 ACM conference on Web science*. ACM, 2014, pp. 42–51.

- [3] L. Hong, O. Dan, and B. D. Davison, "Predicting popular messages in twitter," in *Proceedings of the 20th international conference companion on World wide web*. ACM, 2011, pp. 57–58.
- [4] G. Szabo and B. A. Huberman, "Predicting the popularity of online content," *Communications of the ACM*, vol. 53, no. 8, pp. 80–88, 2010.
- [5] B. Wu, W.-H. Cheng, Y. Zhang, and T. Mei, "Time matters: Multi-scale temporalization of social media popularity," in *Proceedings of the 2016 ACM on Multimedia Conference*. ACM, 2016, pp. 1336–1344.
- [6] Y. Bae and H. Lee, "Sentiment analysis of twitter audiences: Measuring the positive or negative influence of popular twitterers," *Journal of the American Society for Information Science and Technology*, vol. 63, no. 12, pp. 2521–2535, 2012.
- [7] L. De Vries, S. Gensler, and P. S. Leeflang, "Popularity of brand posts on brand fan pages: An investigation of the effects of social media marketing," *Journal of interactive marketing*, vol. 26, no. 2, pp. 83–91, 2012.
- [8] J. Lu, J. Yang, D. Batra, and D. Parikh, "Hierarchical question-image co-attention for visual question answering," in *Advances In Neural Information Processing Systems*, 2016, pp. 289–297.
- [9] Q. You, S. Bhatia, T. Sun, and J. Luo, "The eyes of the beholder: Gender prediction using images posted in online social networks," in *Data Mining Workshop (ICDMW), 2014 IEEE International Conference on*. IEEE, 2014, pp. 1026–1030.
- [10] A. Dhir, S. Pallesen, T. Torsheim, and C. S. Andreassen, "Do age and gender differences exist in selfie-related behaviours?" *Computers in Human Behavior*, vol. 63, pp. 549–555, 2016.
- [11] L. Liu, D. Preotiuc-Pietro, Z. R. Samani, M. E. Moghaddam, and L. H. Ungar, "Analyzing personality through social media profile picture choice," in *ICWSM*, 2016, pp. 211–220.
- [12] Q. You, S. Bhatia, and J. Luo, "A picture tells a thousand words about you! user interest profiling from user generated visual content," *Signal Processing*, vol. 124, pp. 45–53, 2016.
- [13] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent dirichlet allocation," *Journal of machine Learning research*, vol. 3, no. Jan, pp. 993–1022, 2003.
- [14] T. Hu, H. Xiao, J. Luo, and T.-v. T. Nguyen, "What the language you tweet says about your occupation," in *ICWSM*, 2016, pp. 181–190.
- [15] S. Hamidian and M. Diab, "Rumor detection and classification for twitter data," in *Proceedings of the Fifth International Conference on Social Media Technologies, Communication, and Informatics (SOTICS)*, 2015, pp. 71–77.
- [16] N. Kalchbrenner and P. Blunsom, "Recurrent continuous translation models," in *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, 2013, pp. 1700–1709.
- [17] 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 *International Conference on Machine Learning*, 2015, pp. 2048–2057.
- [18] V. Mnih, N. Heess, A. Graves *et al.*, "Recurrent models of visual attention," in *Advances in neural information processing systems*, 2014, pp. 2204–2212.
- [19] M. Qi, J. Qin, A. Li, Y. Wang, J. Luo, and L. Van Gool, "stagnet: An attentive semantic rnn for group activity recognition," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 101–117.
- [20] C. Xiong, V. Zhong, and R. Socher, "Dynamic co-attention networks for question answering," *arXiv preprint arXiv:1611.01604*, 2016.
- [21] J.-H. Kim, S.-W. Lee, D. Kwak, M.-O. Heo, J. Kim, J.-W. Ha, and B.-T. Zhang, "Multimodal residual learning for visual qa," in *Advances in Neural Information Processing Systems*, 2016, pp. 361–369.
- [22] X. Niu, L. Li, T. Mei, J. Shen, and K. Xu, "Predicting image popularity in an incomplete social media community by a weighted bi-partite graph," in *Multimedia and Expo (ICME), 2012 IEEE International Conference on*. IEEE, 2012, pp. 735–740.
- [23] P. J. McParlane, Y. Moshfeghi, and J. M. Jose, "Nobody comes here anymore, it's too crowded; predicting image popularity on flickr," in *Proceedings of International Conference on Multimedia Retrieval*. ACM, 2014, p. 385.
- [24] B. Wu, T. Mei, W.-H. Cheng, Y. Zhang *et al.*, "Unfolding temporal dynamics: Predicting social media popularity using multi-scale temporal decomposition," in *AAAI*, 2016, pp. 272–278.
- [25] M. Mazloom, R. Rietveld, S. Rudinac, M. Worring, and W. Van Dolen, "Multimodal popularity prediction of brand-related social media posts," in *Proceedings of the 2016 ACM on Multimedia Conference*. ACM, 2016, pp. 197–201.
- [26] M. Qi, Y. Wang, and A. Li, "Online cross-modal scene retrieval by binary representation and semantic graph," in *Proceedings of the 2017 ACM on Multimedia Conference*. ACM, 2017, pp. 744–752.
- [27] A. Khosla, A. Das Sarma, and R. Hamid, "What makes an image popular?" in *Proceedings of the 23rd international conference on World wide web*. ACM, 2014, pp. 867–876.
- [28] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [29] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [30] D.-K. Nguyen and T. Okatani, "Improved fusion of visual and language representations by dense symmetric co-attention for visual question answering," *arXiv preprint arXiv:1804.00775*, 2018.

- [31] Q. You, J. Luo, H. Jin, and J. Yang, “Cross-modality consistent regression for joint visual-textual sentiment analysis of social multimedia,” in *Proceedings of the Ninth ACM International Conference on Web Search and Data Mining*. ACM, 2016, pp. 13–22.
- [32] H. Fang, S. Gupta, F. Iandola, R. Srivastava, L. Deng, P. Dollár, J. Gao, X. He, M. Mitchell, J. Platt *et al.*, “From captions to visual concepts and back,” 2015.
- [33] M. Mazloom, B. Hendriks, and M. Worring, “Multimodal context-aware recommender for post popularity prediction in social media,” in *Proceedings of the on Thematic Workshops of ACM Multimedia 2017*. ACM, 2017, pp. 236–244.