Learning User Embeddings from Human Gaze for Personalised Saliency Prediction

FLORIAN STROHM, University of Stuttgart, Germany MIHAI BÂCE*, KU Leuven, Belgium ANDREAS BULLING, University of Stuttgart, Germany

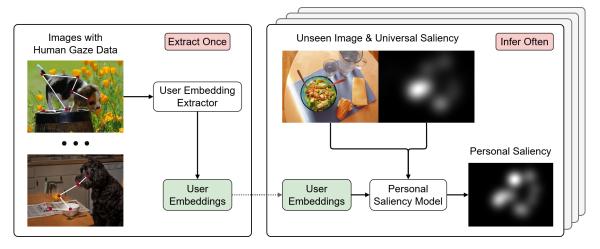


Fig. 1. Our method (left) takes as input a small set of images and their corresponding user-specific saliency maps obtained from human gaze recorded using a stationary eye tracker and produces one user embedding for each user, which captures the user-specific differences in viewing behaviour. We then demonstrate (right) how the learned embeddings can be used to refine a universal saliency map (e.g. obtained from a saliency predictor) to an individual, personal saliency map for any new images.

Reusable embeddings of user behaviour have shown significant performance improvements for the personalised saliency prediction task. However, prior works require explicit user characteristics and preferences as input, which are often difficult to obtain. We present a novel method to extract user embeddings from pairs of natural images and corresponding saliency maps generated from a small amount of user-specific eye tracking data. At the core of our method is a Siamese convolutional neural encoder that learns the user embeddings by contrasting the image and personal saliency map pairs of different users. Evaluations on two public saliency datasets show that the generated embeddings have high discriminative power, are effective at refining universal saliency maps to the individual users, and generalise well across users and images. Finally, based on our model's ability to encode individual user characteristics, our work points towards other applications that can benefit from reusable embeddings of gaze behaviour.

CCS Concepts: • Human-centered computing → User models; • Computing methodologies → Interest point and salient region detections.

Additional Key Words and Phrases: gaze, eye-tracking, saliency, personal saliency, user embeddings, user model, deep learning

1 INTRODUCTION

Saliency prediction is the task of identifying salient regions within an image which are likely to attract gaze. Various models have been developed which take into account both low-level features [26, 27, 57] and high-level image characteristics [16, 37, 43], incorporating bottom-up attention mechanisms, as well as task demands [11, 41, 61], which involve

1

^{*}work conducted while at University of Stuttgart

top-down attention processes. Given the potential for anticipating user attention, saliency prediction models have had significant impact in computer vision and beyond, and have proven highly beneficial for a wide range of tasks, from serving as an inductive bias for neural attention mechanisms [19, 50, 51] to estimating users' cognitive states [1, 25, 56], or enabling personalised predictions for various human-computer-interaction tasks [3, 20, 39, 55, 64].

A large body of work on saliency modelling has focused on *universal saliency*, i.e. the task of predicting saliency maps that aggregate gaze data from multiple observers and, as such, disregard individual differences in viewing behaviour. There are, however, significant individual variations in how visual attention is deployed on image stimuli that are due to a range of factors, such as scene complexity and semantics [13, 62], level of expertise [6, 9, 14, 49], age [63], or personality traits [4, 45].

Despite these differences among individuals and the many applications (e.g. assistive systems [47, 53, 54]) that could benefit from a better understanding of the individual, only few previous works have proposed methods to predict *personalised saliency*. One approach involved training an individual (sub-) model for each user, which lacks generalisability [42, 59, 60]. Another one leverages person-specific information such as age, gender, or preference towards specific object categories or colours [59]. However, in addition to raising privacy concerns, collecting these characteristics is tedious and requires explicit user input.

In contrast to explicitly collecting user information, we introduce a novel method to extract embeddings from users' gaze behaviour while viewing natural images. Our method uses a Siamese convolutional neural encoder that takes multiple images and their corresponding saliency maps of a particular user as input and produces a user embedding as output. The embedding is learned by contrasting input pairs from one user to other users exhibiting different gaze behaviours on the same image stimuli. By integrating the user embedding into a saliency prediction network (Figure 1), this additional input plays a crucial role in predicting filters that are convolved over extracted image features, thus integrating user-specific information essential for the personalised saliency prediction task. Our findings demonstrate the highly discriminative nature of these embeddings, enabling us to effectively compare individuals based on their distinctive gaze behaviour (Figure 6). Results on the downstream task of personal saliency prediction task show how the generated embeddings can be used to effectively refine the universal saliency map predictions and tailor them to individual users. Moreover, we observe that these embeddings exhibit good generalisation capabilities when applied to both unseen users and images. ¹

2 RELATED WORK

2.1 Individual Differences in Visual Saliency

Traditional saliency prediction methods ignore individual differences in visual salience between humans and instead predict an average, universal saliency map [5]. However, there are multiple prior works that show that humans have different visual preferences which draws their attention, which are stable and predictable. De Haas and Linka et al. [13, 38] have identified multiple semantic dimensions along which human salience significantly differs, like faces. Prior works have incorporated face detectors in saliency prediction pipelines, as they generally tend to attract significant attention [5, 10]. However, the results from De Haas et al. show that for specific humans, such predictions are imprecise as their attention is not attracted by faces. Later Broda et al. [7, 8] identified that humans can be roughly clustered in two categories when observing persons. Either they tend to focus the head and inner facial features or they fixate on body parts like arms and legs. Prior work has also studied which human traits influence the individual attention and

¹Project code will be made publicly available upon acceptance.

found that for example age [33] is an important factor as well as personality [23] and gender [46]. Xu et al. [59] were the first that proposed a method to predict personalised instead of universal saliency by utilising such user traits as an additional input to their network. However, it is still unclear which user traits are useful for saliency prediction, and explicitly collecting personal information might not be appropriate. Moroto et al. [42] later proposed a method involving a multi-task CNN to predict personalised saliency for each user in the training dataset. During inference, unseen users were matched to the seen users based on their similarity in attention allocation. However, this requires many different users in the training data to generalise to unseen users. Moreover, finding an appropriate similarity function to match unseen to seen users based on their gaze behaviour is challenging. In contrast, we propose to train a single-task CNN that incorporates an additional user embedding as input, enabling the network to leverage user-specific information.

2.2 User Embeddings

User embeddings have found applications across diverse domains, serving either as a means to infer user-related information or to personalise the user experience. In a seminal work by Pazzani et al. [44], an approach was introduced where an agent learns a user embedding by leveraging explicit feedback provided by users in the form of page ratings. This user embedding was then utilised to provide personalised website recommendations tailored to the user's interests and preferences. Later, various methods have been proposed to construct user embeddings by incorporating diverse forms of application-specific explicit feedback [34, 48]. Methods that collect user information explicitly require the user to be active and might become inaccurate over time as the user's interest changes.

To address this limitation, implicit methods to create user embeddings have gained popularity. These methods allow users to simply interact with a system while the system creates a user embedding, without relying on explicit feedback [15, 17]. In a recent work by Wu et al. [58], a novel approach called Author2Vec was introduced to derive user embeddings by analysing the textual content authored by users on social media platforms. The authors demonstrated the superior performance of these embeddings in predicting user personality traits or mental states compared to alternative methods. Similarly, An et al. [2] used web browsing events to extract rich user embeddings for different downstream tasks. In other related works, researchers investigated the incorporation of user embeddings as unique word tokens alongside the user's text. This approach enables the model to comprehend the sentence in the context of the user embedding [40, 65]. In the eye-gaze domain He et al. [18] have utilised appearance based user embeddings extracted from faces to personalise an appearance-based gaze estimator. These works show the potential to leverage implicit user embeddings on a range of applications. In our work, to the best of our knowledge, we are the first to learn user embeddings from visual attentive behaviour and leverage them to enhance performance on the personalised saliency prediction task.

3 METHODOLOGY

Traditional saliency prediction models learn a function $f(I) = \text{USM}_I$, where I is an image and USM_I is the corresponding Universal Saliency Map (USM). The ground truth USM is calculated by averaging n Personal Saliency Maps (PSM) obtained from different humans observing the same image: $\text{USM}_I = \frac{1}{n} \sum_{k=1}^n \text{PSM}_{I,k}$. Instead of predicting USM_I , our goal is to predict $\text{PSM}_{I,U}$ for a given image I and user U. To predict $\text{PSM}_{I,U}$ it is essential to incorporate additional input information that is specific to the user U for whom the predictions are intended.

3

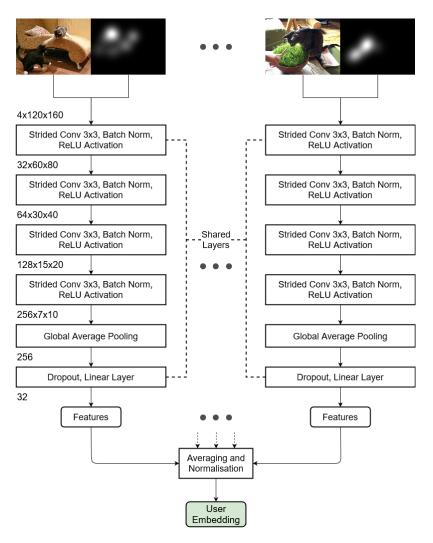


Fig. 2. The architecture of our proposed user embedding extractor involves processing multiple images alongside an additional channel that includes the saliency information of a specific user, from which we aim to extract an embedding. To accomplish this, we employ a Siamese convolutional neural network, which is responsible for extracting features from each pair of image-saliency maps. Subsequently, the extracted features are averaged and normalised, resulting in the user embedding.

3.1 Learning User Embeddings

We propose a novel method for extracting user embeddings e_U for each individual user U. These embeddings are derived from the distinct variations observed in the visual attention patterns of individual users. Our hypothesis is that these embeddings can be effectively used by a personalised saliency prediction model to generate user-specific saliency outputs.

The architecture of our proposed embedding neural network is shown in Figure 2. A Siamese user embedding extractor E takes m different images $\{I_1, ..., I_m\}$ along with the corresponding PSMs $\{PSM_{I_1,U}, ..., PSM_{I_m,U}\}$ for the same user U as input. The goal is to extract joint image-saliency features which allow the network to understand

the visual preferences of the user and extract a meaningful user embedding. The PSM for each image is treated as an additional image channel besides the existing three RGB channels and resized to a resolution of 160×120 . Each image-PSM tensor is passed through four convolution blocks consisting of a 2D convolution layer with stride two, a batch normalisation layer [24] and a Rectified Linear Unit (ReLU) activation function. Subsequently, a global average pooling layer [36] is employed to reduce the output to a one-dimensional vector, which helps prevent overfitting in conjunction with a dropout layer [52]. Finally, a linear layer predicts the output features for each image-PSM input.

This Siamese network is used to extract joint image-saliency features for each of the *m* image-PSM pairs, which are subsequently averaged and normalised to unit length resulting in the extracted user embedding. Preliminary experiments revealed that combining the extracted features from each image-PSM pair with a recurrent or transformer network results in severe overfitting. Thus, averaging the features helps prevent overfitting and yields better results overall.

The network is optimised to minimise the triplet margin loss with online (semi-) hard triplet mining [21] defined as:

$$e_{a} = \|E(\{I_{1}, ..., I_{m}\}, \{PSM_{I_{1}, U_{1}}, ..., PSM_{I_{m}, U_{1}}\})\|$$

$$e_{p} = \|E(\{I_{m+1}, ..., I_{2m}\}, \{PSM_{I_{m+1}, U_{1}}, ..., PSM_{I_{2m}, U_{1}}\})\|$$

$$e_{n} = \|E(\{I_{2m+1}, ..., I_{3m}\}, \{PSM_{I_{2m+1}, U_{2}}, ..., PSM_{I_{3m}, U_{2}}\})\|$$

$$\mathcal{L}_{E} = \max(e_{a} \cdot e_{p} - e_{a} \cdot e_{n} + m, 0).$$

$$(1)$$

To calculate the anchor user embedding e_a in Equation (1), a random user U_1 and m random images $\{I_1,...,I_m\}$ with corresponding PSMs $\{PSM_{I_1,U_1},...,PSM_{I_m,U_1}\}$ are selected and passed through our embedding network. Similarly, the positive embedding example e_p is calculated using the same user U_1 but with different images $\{I_{m+1},...,I_{2m}\}$ and PSMs $\{PSM_{I_{m+1},U_1},...,PSM_{I_{2m},U_1}\}$, while the negative embedding example e_n is obtained by selecting a different random user U_2 , also with different images $\{I_{2m+1},...,I_{3m}\}$ and corresponding PSMs $\{PSM_{I_{2m+1},U_2},...,PSM_{I_{3m},U_2}\}$. Based on the three embeddings e_a , e_p and e_n the standard triplet loss can be calculated as defined Equation (1). It is crucial to note that each image I in the training dataset has a corresponding personalised saliency map $PSM_{I,U}$ for each user U. This ensures that the user embeddings are solely derived from the individual differences in users' visual attention behaviour, rather than being influenced by the images themselves. By minimising the distance between the anchor and the positive embedding, the network learns to recognise similar high-level gaze behaviour between different inputs for the same user. Similarly, maximising the distance between the anchor and the negative embedding encourages the network to distinguish the different attentive behaviour between two users.

While the training objective is to differentiate between different users, prior research has shown that embeddings learned through optimising the triplet loss have the capability to encode substantial class-specific information [21]. Additionally, similar classes tend to be close to each other in the embedding space, while dissimilar classes tend to have greater separation. This indicates that the embedding space effectively captures the relevant characteristics of the classes and presents a structured representation that reflects the underlying relationships between them. By using the user embedding, a downstream task model can leverage the captured information and tailor its predictions to the specific characteristics and visual attention behaviour of each user.

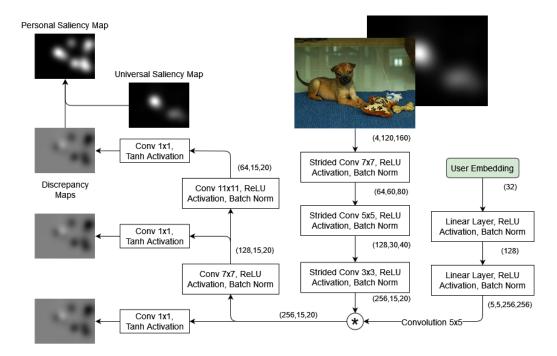


Fig. 3. The personalised saliency map (PSM) network operates by taking an image stimulus and its corresponding universal saliency map (USM) as input. In addition, it incorporates user embedding, which is utilised to predict kernel weights, which are then convolved over the image-USM features. The network outputs a discrepancy map which can be added to the USM in order to generate the PSM.

3.2 Personalised Saliency Prediction

To demonstrate the effectiveness of the user embeddings, we adapted an existing method for personalised saliency prediction [59] and replace the manually, explicitly defined user characteristics with our learned embeddings from implicit gaze behaviour.

Similar to Xu et al [59] we define the PSM prediction task as a refinement of the USM:

$$PSM(I, U) = USM(I) + \Delta(I, U), \tag{2}$$

where PSM(I,U) is the PSM of a user U for a given image I, USM(I) is the USM for that image and $\Delta(I,U)$ is the discrepancy map of that user and that image. The discrepancy map essentially defines which parts of the image attract the specific users attention more or less compared to the average user. Using this definition allows us to disentangle the prediction of the PSM, focusing on predicting the discrepancy map $\Delta(I,U)$ while using existing state-of-the-art models to predict the USM. Figure 3 visualises our neural network architecture for personalised saliency prediction. Input to the network N are the image, the USM as a forth image channel and a user embedding. First, a series of convolution layers extract image-saliency features. In parallel, a series of linear layers predict filters based on the user embedding, which are subsequently convolved over the extracted image-saliency features. This way the network can learn features that are specifically important to the user based on that users embedding [59]. After a final sequence of convolution layers the model outputs a discrepancy map of size 15×20 with a value range between -1 and 1. We optimise the network with the mean squared error (MSE) loss between the predicted and the target discrepancy map. Previous work has

shown that by projecting extracted features from intermediate layers to the output space for additional supervision can improve performance [12, 35]. Therefore we add an additional projection layer before the last two convolution layers to predict intermediate discrepancy maps. Our overall loss objective \mathcal{L}_{PSM} is then given as:

$$\mathcal{L}_{PSM} = \left\| \sum_{k=1}^{3} N(\text{USM}(I), I, E_U)_k - \Delta(I, U) \right\|,\tag{3}$$

where $N(...)_k$ is one of the three predicted discrepancy maps of the personalised saliency network. To calculate the ground truth discrepancy maps $\Delta(I, U)$ we use Equation (2) and subtract the USM_I from the PSM_{I,U}.

4 EXPERIMENTS

4.1 Implementation Details

User Embedding Network. The architecture of the user embedding network is shown in Figure 2. The number of input image-PSM samples m to our proposed Siamese CNN depends on the experiment and varies between 4 and 32. Each model was optimised with the triplet margin loss as defined in Equation (1) with a margin of 0.05 and online semi-hard and hard triplet mining within a batch of 256 samples. The weights were updated using the Adam optimiser [32] with a learning rate of 0.001 and default parameters otherwise. The Dropout layer [52] masked neuron activations 50% of the time during training.

Personalised Saliency Network. Figure 3 illustrates the architecture of our personalised saliency network. To effectively train this network, we first trained the embedding network to extract user embeddings, which serve as input. Since the embeddings rely on *m* image-PSM pairs for their calculation, they can exhibit variability, particularly when *m* is small. Hence, after training the embedding networks, we generated 100 embeddings for each participant within each dataset by randomly sampling data from the corresponding participant. During the training process of the personalised saliency network, we randomly selected embeddings from this pool to account for these variations and improve generalisation. The models were trained with an SGD optimiser with initial learning rate of 0.02, momentum of 0.9, weight decay of 0.0005 and batch size 32. The learning rate was reduced by a factor of 2 every 25 epochs. The trained model runs at 27 frames per second on an RTX 4070 GPU, achieving real-time personalised saliency predictions.

Evaluation Metrics. We report multiple metrics commonly used to evaluate the similarity between saliency maps [30, 43, 59]. These metrics are the Pearson's correlation coefficient (CC), similarity / histogram intersection (SIM), Area under ROC Curve (AUC-Judd) [31], normalised scanpath saliency (NSS) and Kullback-Leibler divergence (KLD). In addition, to assess the accuracy of the user embeddings, we employ a labelling approach where an extracted embedding is considered correct if its nearest neighbour in the embedding space belongs to the same user. This evaluation metric is commonly referred to as precision at one [28].

4.2 Datasets

Saliency prediction models are typically pre-trained using the large-scale SALICON dataset [29]. However, we cannot utilise the SALICON dataset for learning the embeddings as each subject from the dataset observed a different subset of image stimuli. As discussed in Subsection 3.1, it is critical that each participant looks at all images, or at least that there is a significant overlap between participants, as otherwise the model can simply identify the user based on the images they observed ignoring their specific visual attention behaviour.

| Model | CC | SIM | AUC | NSS | KLD |
|------------------------------------|-----------------------|-----------------------|-----------------------|----------------|-----------------------|
| DeepGazeIIE (DG) | 0.622 | 0.556 | 0.904 | 2.121 | 1.123 |
| Fine-tuned DG | 0.715 | 0.619 | 0.905 | 2.140 | 0.526 |
| MultiCNN w/ DG | 0.735 | 0.643 | 0.897 | 2.142 | 0.708 |
| Ours w/ DG | 0.736 | 0.651* | 0.907 | 2.170* | 0.509* |
| | | | | | |
| Ground Truth USM | 0.801 | 0.685 | 0.921 | 2.373 | 0.372 |
| Ground Truth USM MultiCNN w/ GT | 0.801 <u>0.804</u> | $\frac{0.685}{0.683}$ | $\frac{0.921}{0.919}$ | 2.373 2.378 | $\frac{0.372}{0.511}$ |

Table 1. Closed-set results for the ID dataset [13]. * indicates significant improvement over the strongest baseline.

| Model | CC | SIM | AUC | NSS | KLD |
|------------------|--------|--------------|--------------|-------|--------------------|
| DeepGazeIIE (DG) | 0.584 | 0.530 | 0.881 | 2.288 | 1.294 |
| Fine-tuned DG | 0.734 | 0.622 | 0.887 | 2.293 | 0.562 |
| MultiCNN w/ DG | 0.746 | 0.637 | 0.892 | 2.299 | $\overline{0.604}$ |
| Ours w/ DG | 0.760* | 0.649* | 0.896 | 2.308 | 0.481* |
| Ground Truth USM | 0.821 | 0.698 | 0.912 | 2.500 | 0.366 |
| | | | | | |
| MultiCNN w/ GT | 0.845 | <u>0.711</u> | <u>0.914</u> | 2.609 | 0.415 |

Table 2. Closed-set results for the PS dataset [59]. * indicates significant improvement over strongest baseline.

Based on the above requirement, we selected two publicly available datasets where each participant observed each image while their gaze was recorded with an eye-tracker. The first dataset was collected by Xu et al. [59], which we call the *Personalised Saliency* (PS) dataset. The PS dataset contains 1,600 images, which they selected to contain many different semantic categories in each image, as they argue that this maximises the variation of visual attention between participants. Each stimuli in the dataset was observed by 30 participants for three seconds a total of four times, allowing them to average out stochastic variations in each participants attentive behaviour. To evaluate how well our system generalises to unseen images we split the images into 80% for training, 10% for validation and 10% for testing. Furthermore, to evaluate how well our system generalises to unseen participants, we split the 30 participants into 20 for training and 5 each for validation and testing.

The second dataset was collected by Haas et al. [13] which we call the *Individual Differences* (ID) dataset. The stimuli were selected to be comprised of complex scenes containing multiple different semantic categories and objects. They recorded gaze data from a total of 102 different participants each looking at 700 different image stimuli. Similar to the PS dataset we split the images into 80% for training, 10% for validation and 10% for testing. Furthermore, we split the 102 participants into 80 for training, 10 for validation and 12 testing.

Following Xu et al. [59] we conduct both closed-set and open-set experiments for each dataset. In the closed-set experiments, the same participants were used in both the training, validation, and test sets with the validation and test sets containing unseen images. This enables us to assess the capability to predict personalised saliency for familiar users. As for the open-set experiments, we assess the personalised saliency prediction performance on unseen participants from the test set. This analysis helps us understand how effectively models can generalise to new users.

| Model | CC | SIM | AUC | NSS | KLD |
|-----------------------------|--------------------|-----------------------|-----------------------|-------|--------------------------|
| DeepGazeIIE (DG) | 0.615 | 0.553 | 0.916 | 2.239 | 1.083 |
| Fine-tuned DG | 0.721 | 0.624 | 0.914 | 2.244 | 0.542 |
| MultiCNN w/ DG | 0.723 | 0.627 | 0.910 | 2.244 | 0.620 |
| | | | | | |
| Ours w/ DG | 0.726 | 0.638* | 0.916 | 2.253 | 0.527^* |
| Ours w/ DG Ground Truth USM | 0.726 0.804 | 0.638 * 0.682 | 0.916 0.922 | 2.253 | |
| · | | | | | 0.527* 0.387 0.440 |

Table 3. Open-set results for the ID dataset [13]. * indicates significant improvement over strongest baseline.

| Model | CC | SIM | AUC | NSS | KLD |
|------------------------------------|-------------------------|-------------------------|-------------------------|-----------------------|-------------------------|
| DeepGazeIIE (DG) | 0.564 | 0.522 | 0.860 | 1.944 | 1.431 |
| Fine-tuned DG | 0.692 | 0.618 | 0.862 | 2.032 | 0.685 |
| MultiCNN w/ DG | 0.693 | 0.615 | 0.862 | 2.006 | 0.682 |
| Ours w/ DG Ours CD w/ DG | 0.712 0.713 * | 0.634 0.634 * | 0.870 0.871 * | 2.033 2.033 | 0.539 0.533 * |
| Ground Truth USM MultiCNN w/ GT | 0.788 <u>0.800</u> | 0.691 0.694 | 0.892 0.892 | 2.218 2.290 | 0.399 0.516 |
| Ours w/ GT Ours CD w/ GT | 0.804 0.807 | 0.701 0.704 * | 0.893 0.894 | 2.273 2.288 | 0.371 0.361 * |

Table 4. Open-set results for the PS dataset [59]. * indicates significant improvement over strongest baseline.

4.3 Baselines

We evaluate our method against three baselines. Firstly, we utilise DeepGaze IIE [37], a state-of-the-art universal saliency prediction model. We compare the predicted USMs by DeepGaze IIE with our refined PSMs based on the DeepGaze IIE prediction. Since DeepGaze IIE was not originally trained on our dataset, we fine-tune the model's prediction on each dataset and present results both with and without fine-tuning. Secondly, we compare the ground truth USMs with our predicted PSMs generated through refining the USMs.

In the literature, there are two relevant works by Xu et al. [59] and Moroto et al. [42] that propose methods for PSM prediction, as discussed in Section 2. Unfortunately, we were unable to directly compare our results with Xu et al. due to the unavailability of the user-specific information required for their method. However, for the closed set experiments they also proposed MultiCNN, where they train a separate classifier for each participant, which we will use as our third baseline. The MultiCNN architecture is identical to our network shown in Figure 3 without the embedding pathway.

Moroto et al. [42] propose a method to map unseen users to the training user with the most similar visual attention behaviour, which allows them to make personalised saliency predictions with the corresponding trained classifier. Inspired by this, we propose an oracle mapping by evaluating unseen participants using every trained MultiCNN model and then choose the model achieving the lowest loss. This allows us to provide the upper-bound MultiCNN performance for the open set experiments.

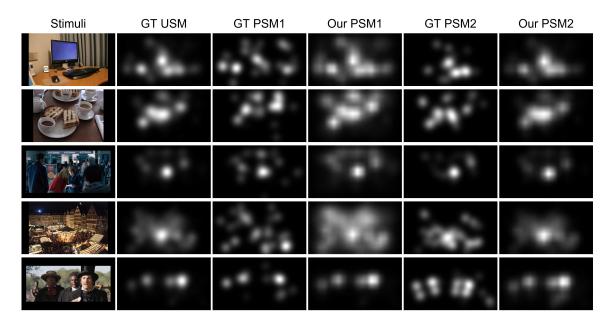


Fig. 4. Example PSM predictions for two users from the PS [59] test set with our proposed method compared to the ground truths.

4.4 Personalised Saliency Prediction

Closed-Set Results. Table 1 shows the closed-set results for the ID dataset and Table 2 shows the results for the PS dataset. The results with our method were obtained using embeddings extracted from m=32 image-PSM pairs for each user. The best performing models are highlighted in bold, while the performance of the second best model is underlined. Significance tests were conducted using the Mann-Whitney-U test for each metric with a p-value threshold of < 0.05. An asterisk next to metrics for our method indicates that the improvement compared to the strongest baseline was statistically significant. The ground truth USM is the upper bound traditional saliency prediction methods could potentially achieve. The results for MultiCNN and Ours show that it is possible to further refine these USMs as overall both methods outperform the ground truth USM baseline. Furthermore, we observe that for both datasets our method outperforms all baselines including MultiCNNs in all metrics except for NSS on the PS dataset. However, for a real-world scenario the ground truth USM might not be available and has to be predicted first. The results show that MultiCNNs and Ours using the non fine-tuned USM predictions from DeepGaze IIE still outperform the fine-tuned DeepGaze IIE baseline, with our method achieving the best performance. This indicates that our method can be applied to different potential suboptimal USM predictions and still produce a more refined PSM prediction.

Open-Set Results. Table 3 shows the open-set results for the ID dataset and Table 4 shows the results for the PS dataset. We can observe a very similar trend as with the closed-set experiments with our method overall outperforming all baselines, indicating that our embeddings help the personal saliency network to generalise well to unseen participants. Since the PS dataset consists only of 30 different participants of which 20 are used for training, learning generalisable user embeddings is more challenging. We therefore experiment with combining the training splits of the PS and ID datasets when training the user embedding network, allowing the network to observe the attentive behaviour of a total of 100 different participants during training. Note that we still trained and evaluated our personalised saliency network

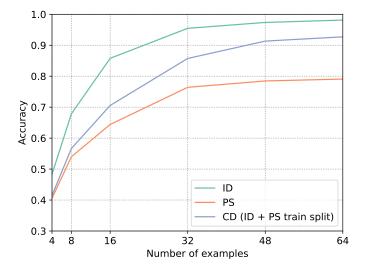


Fig. 5. Performance comparison of different user embedding extraction models. The y-axis indicates the model's accuracy and the x-axis how many image-PSM examples m were used as input to extract the embedding (4, 8, 16, 32, 48 or 64). We report the accuracy for unseen participants on the Individual Differences (ID) dataset, the Personal Saliency (PS) dataset and for the combined CD dataset.

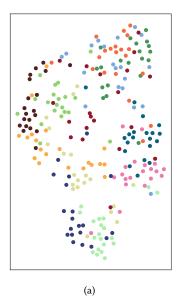
using only the PS dataset. We report the resulting performance on the combined dataset (CD) as *Ours CD* in Table 4. We can observe that using these new embeddings the performance of the personalised saliency model further improves for all metrics.

In addition to the quantitative results, Figure 4 shows example open-set PSM predictions for multiple image stimuli on the test split of the PS dataset. In the provided examples, it is evident that the user represented in the last two columns exhibits a more focused attention behaviour compared to the first user. Our embeddings successfully capture this distinction, as indicated by the corresponding saliency predictions. Notably, in the example depicted in the bottom row, the second user appears to allocate less attention to faces and instead also focuses on other body parts. Our embeddings seem to capture this behaviour, as evidenced by the saliency allocation below the face region in the corresponding PSM predictions. Together with our quantitative results this further demonstrates the effectiveness of using user embeddings learned from visual attention for more personalised saliency predictions.

4.5 User Embeddings Analysis

To gain a better understanding of our extracted user embeddings, we further analyse the performance of the user embedding network.

Number of Examples for Embedding Extraction. Figure 5 shows the model's test set accuracy on the y-axis for different datasets and the number of examples m used for embedding extraction on the x-axis. We report the test set accuracy for both datasets ID and PS, as well as the PS test set accuracy when training on the combined training set CD. The user embedding network achieves an accuracy of 98.1% on the ID test split when using m = 64 examples, showing that it is able to differentiate very well between participants that the model never saw during training. Similarly the model trained on the PS dataset achieves an accuracy of 79.1% when only using the PS training split and an accuracy of



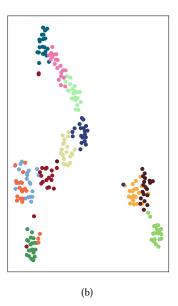


Fig. 6. We use t-SNE to reduce the 32-dimensional embeddings into two dimensions. Figure (a) shows the embedding space for embeddings extracted using m = 8 image-PSM pairs while in (b) we used m = 32 pairs. Each dot represents one embedding extracted using m randomly sampled image-PSM pairs. Each colour corresponds to a unique user from the ID test set.

92.7% when combining the training splits from PS and ID. As already indicated by the improved personalised saliency prediction results reported in Table 4, training on both datasets increases the performance on the PS dataset, which only contains a small number of participants. Note that the random baseline accuracy for the ID dataset is 8.3% (12 participants in the test set), while it is 20% for the PS dataset (5 participants in the test set). Furthermore, we can observe that the model's accuracy increases with the number of examples m provided as input. Combining the training datasets especially improves the accuracy for larger m. This is likely due to the increased user information the model can extract with larger m, resulting in better generalisation if provided with more diverse users during training.

Visualisation of the Embedding Space. To analyse the embedding space we reduced the 32 dimensional user embeddings to 2 dimensions using t-SNE. Figure 6 illustrates the 2-dimensional embeddings for the 12 participants from the ID test split. In this visualisation, each dot represents an embedding, and the colour of each dot corresponds to a specific user. We calculated 20 embeddings for each participant by randomly sampling m input image-PSM pairs. In (a) we visualise the embeddings for m = 8 while (b) shows the embeddings for m = 32. We observe that the embeddings for m = 8 have a much larger variance within a participant compared to the embeddings extracted with m = 32, as less information about the user can be extracted. For a low m the boundaries between participants becomes fuzzy which is also reflected by the lower accuracy reported in Figure 5. We can observe that for m = 32 the embeddings are highly discriminative as clear user clusters are formed.

5 BROADER IMPACT

Our method demonstrates the ability to extract user embeddings from a small amount of gaze data, showcasing that these embeddings effectively capture relevant information for modelling personal saliency. Predicting saliency is a

crucial task with implications for various downstream applications in computer vision and human-computer interaction. For instance, attentive user interfaces aim to manage the user's attention effectively, relying on knowledge about their visual attention [55]. Recommender systems can benefit from understanding users' visual attention to enhance the visibility of top-ranked entities [64]. Other downstream tasks that benefit from accurate personalised saliency include video summarization [39], automated image cropping [3], and image captioning [20]. While the downstream task of this work is personalised saliency prediction, our proposed method for extracting user embeddings could prove advantageous for other downstream tasks where personalisation is crucial.

Although personalised saliency is helpful for many tasks, the underlying computational user model could also be misused to synthesise data for a given user by impersonating the user's own saliency. Furthermore, instead of using the embeddings to predict personal saliency, they might be misused to directly extract user sensitive information. These embeddings might encode user characteristics or private information about the user that correlates with their gaze behaviour. For example, prior work has shown that gender is a strong modulator of saliency preference [22], and as such, this private information might be extractable from our embeddings. This encoded information could be used to identify users based on their personal saliency, especially as attention tracking becomes more pervasive and cheap through appearance based gaze estimation or mouse tracking. Another factor to consider is that our embeddings may be biased as they are extracted from a small number of stimuli with corresponding gaze data. As we can see in Figure 6, the embeddings vary for the same person and we can notice multiple outlier embeddings that are far apart form their cluster, even when using more data. Thus, when used for tasks like user profiling this might lead to incorrect user profiles and subsequently incorrect predictions in downstream tasks.

Continuing this line of research, it is conceivable that future work will not only synthesise personalised saliency but even raw gaze of specific users with such embeddings. Accurate computational models of user-specific gaze behaviour would be of great significance for even more downstream tasks. However, this advancement might raise potential security risks for applications that fundamentally rely on gaze behaviour analysis, such as gaze-based authentication.

6 SUMMARY

In this work, we proposed a novel method that extracts user embeddings from pairs of natural images and corresponding user-specific saliency maps. The learned embeddings capture the users' unique characteristics and can be used to address the personalised saliency prediction task. In contrast to prior work for this task that required explicit user input, our method only requires implicit input from gaze behaviour collected using an eye tracker. Our proposed method uses a Siamese convolutional neural encoder to learn the embedding model, trained by contrasting a user's gaze behaviour with that of different users. Results on two saliency datasets demonstrated the embeddings' discriminative power, our method's generalisability to unseen users and images, and improved performance over universal saliency prediction models. As such, our work presents a promising approach to learning and leveraging user embeddings from implicit behaviour also for other tasks or applications that require individual user characteristics.

ACKNOWLEDGMENTS

Florian Strohm and Andreas Bulling were funded by the European Research Council (ERC) under the grant agreement 801708. Mihai Bâce was funded by a Swiss National Science Foundation (SNSF) Postdoc.Mobility Fellowship (grant number 214434).

REFERENCES

- [1] Ahmed Abdou, Ekta Sood, Philipp Müller, and Andreas Bulling. 2022. Gaze-enhanced Crossmodal Embeddings for Emotion Recognition. *Proceedings of the ACM on Human-Computer Interaction* 6, ETRA (2022), 1–18.
- [2] Mingxiao An and Sundong Kim. 2021. Neural user embedding from browsing events. In Machine Learning and Knowledge Discovery in Databases: Applied Data Science Track: European Conference, ECML PKDD 2020, Ghent, Belgium, September 14–18, 2020, Proceedings, Part IV. Springer, 175–191.
- [3] Edoardo Ardizzone, Alessandro Bruno, and Giuseppe Mazzola. 2013. Saliency based image cropping. In Image Analysis and Processing-ICIAP 2013: 17th International Conference, Naples, Italy, September 9-13, 2013. Proceedings, Part I 17. Springer, 773–782.
- [4] Adrien Baranes, Pierre-Yves Oudeyer, and Jacqueline Gottlieb. 2015. Eye movements reveal epistemic curiosity in human observers. Vision research 117 (2015), 81–90.
- [5] Ali Borji, Ming-Ming Cheng, Qibin Hou, Huaizu Jiang, and Jia Li. 2019. Salient object detection: A survey. Computational visual media 5 (2019), 117–150.
- [6] Stephanie Brams, Gal Ziv, Oron Levin, Jochim Spitz, Johan Wagemans, A Mark Williams, and Werner F Helsen. 2019. The relationship between gaze behavior, expertise, and performance: A systematic review. Psychological bulletin 145, 10 (2019), 980.
- [7] Maximilian Davide Broda and Benjamin De Haas. 2022. Individual differences in looking at persons in scenes. Journal of Vision 22, 12 (2022), 9-9.
- [8] Maximilian D Broda and Benjamin de Haas. 2022. Individual fixation tendencies in person viewing generalize from images to videos. *i-Perception* 13, 6 (2022), 20416695221128844.
- [9] Guy Thomas Buswell. 1935. How people look at pictures: a study of the psychology and perception in art. (1935).
- [10] Moran Cerf, Jonathan Harel, Wolfgang Einhäuser, and Christof Koch. 2007. Predicting human gaze using low-level saliency combined with face detection. Advances in neural information processing systems 20 (2007).
- [11] Xianyu Chen, Ming Jiang, and Qi Zhao. 2021. Predicting human scanpaths in visual question answering. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 10876–10885.
- [12] Marcella Cornia, Lorenzo Baraldi, Giuseppe Serra, and Rita Cucchiara. 2016. A deep multi-level network for saliency prediction. In 2016 23rd International Conference on Pattern Recognition (ICPR). IEEE, 3488–3493.
- [13] Benjamin De Haas, Alexios L Iakovidis, D Samuel Schwarzkopf, and Karl R Gegenfurtner. 2019. Individual differences in visual salience vary along semantic dimensions. *Proceedings of the National Academy of Sciences* 116, 24 (2019), 11687–11692.
- [14] Shahram Eivazi, Roman Bednarik, Markku Tukiainen, Mikael von und zu Fraunberg, Ville Leinonen, and Juha E Jääskeläinen. 2012. Gaze behaviour of expert and novice microneurosurgeons differs during observations of tumor removal recordings. In *Proceedings of the Symposium on Eye Tracking Research and Applications*. 377–380.
- [15] Christopher Ifeanyi Eke, Azah Anir Norman, Liyana Shuib, and Henry Friday Nweke. 2019. A survey of user profiling: State-of-the-art, challenges, and solutions. IEEE Access 7 (2019), 144907–144924.
- [16] Camilo Fosco, Anelise Newman, Pat Sukhum, Yun Bin Zhang, Nanxuan Zhao, Aude Oliva, and Zoya Bylinskii. 2020. How much time do you have? modeling multi-duration saliency. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 4473–4482.
- [17] Susan Gauch, Mirco Speretta, Aravind Chandramouli, and Alessandro Micarelli. 2007. User profiles for personalized information access. The adaptive Web: methods and strategies of Web personalization (2007), 54–89.
- [18] Junfeng He, Khoi Pham, Nachiappan Valliappan, Pingmei Xu, Chase Roberts, Dmitry Lagun, and Vidhya Navalpakkam. 2019. On-device few-shot personalization for real-time gaze estimation. In Proceedings of the IEEE/CVF international conference on computer vision workshops. 0–0.
- [19] Shengfeng He, Chu Han, Guoqiang Han, and Jing Qin. 2019. Exploring duality in visual question-driven top-down saliency. *IEEE transactions on neural networks and learning systems* 31, 7 (2019), 2672–2679.
- [20] Sen He, Hamed R Tavakoli, Ali Borji, and Nicolas Pugeault. 2019. Human attention in image captioning: Dataset and analysis. In Proceedings of the IEEE/CVF International Conference on Computer Vision. 8529–8538.
- [21] Alexander Hermans, Lucas Beyer, and Bastian Leibe. 2017. In defense of the triplet loss for person re-identification. arXiv preprint arXiv:1703.07737
- [22] Johannes Hewig, Ralf H Trippe, Holger Hecht, Thomas Straube, and Wolfgang HR Miltner. 2008. Gender differences for specific body regions when looking at men and women. *Journal of Nonverbal Behavior* 32 (2008), 67–78.
- [23] Sabrina Hoppe, Tobias Loetscher, Stephanie A Morey, and Andreas Bulling. 2018. Eye movements during everyday behavior predict personality traits. Frontiers in human neuroscience (2018), 105.
- [24] Sergey Ioffe and Christian Szegedy. 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In International conference on machine learning. pmlr, 448–456.
- [25] Shamsi T Iqbal and Brian P Bailey. 2004. Using eye gaze patterns to identify user tasks. In The Grace Hopper Celebration of Women in Computing, Vol. 4. 2004.
- [26] Laurent Itti and Christof Koch. 2001. Computational modelling of visual attention. Nature reviews neuroscience 2, 3 (2001), 194–203.
- [27] Laurent Itti, Christof Koch, and Ernst Niebur. 1998. A model of saliency-based visual attention for rapid scene analysis. IEEE Transactions on pattern analysis and machine intelligence 20, 11 (1998), 1254–1259.
- [28] Kalervo Järvelin and Jaana Kekäläinen. 2017. IR evaluation methods for retrieving highly relevant documents. In ACM SIGIR Forum, Vol. 51. ACM New York, NY, USA, 243–250.

- [29] Ming Jiang, Shengsheng Huang, Juanyong Duan, and Qi Zhao. 2015. Salicon: Saliency in context. In Proceedings of the IEEE conference on computer vision and pattern recognition. 1072–1080.
- [30] Tilke Judd, Frédo Durand, and Antonio Torralba. 2012. A benchmark of computational models of saliency to predict human fixations. (2012).
- [31] Tilke Judd, Krista Ehinger, Frédo Durand, and Antonio Torralba. 2009. Learning to predict where humans look. In 2009 IEEE 12th international conference on computer vision. IEEE. 2106–2113.
- [32] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014).
- [33] Onkar Krishna, Andrea Helo, Pia Rämä, and Kiyoharu Aizawa. 2018. Gaze distribution analysis and saliency prediction across age groups. PloS one 13, 2 (2018), e0193149.
- [34] Bruce Krulwich. 1997. Lifestyle finder: Intelligent user profiling using large-scale demographic data. Al magazine 18, 2 (1997), 37-37.
- [35] Chen-Yu Lee, Saining Xie, Patrick Gallagher, Zhengyou Zhang, and Zhuowen Tu. 2015. Deeply-supervised nets. In Artificial intelligence and statistics. PMLR. 562–570.
- [36] Min Lin, Qiang Chen, and Shuicheng Yan. 2013. Network in network. arXiv preprint arXiv:1312.4400 (2013).
- [37] Akis Linardos, Matthias Kümmerer, Ori Press, and Matthias Bethge. 2021. DeepGaze IIE: Calibrated prediction in and out-of-domain for state-of-the-art saliency modeling. In Proceedings of the IEEE/CVF International Conference on Computer Vision. 12919–12928.
- [38] Marcel Linka and Benjamin de Haas. 2020. OSIEshort: A small stimulus set can reliably estimate individual differences in semantic salience. Journal of vision 20, 9 (2020), 13–13.
- [39] Yu-Fei Ma, Lie Lu, Hong-Jiang Zhang, and Mingjing Li. 2002. A user attention model for video summarization. In *Proceedings of the tenth ACM international conference on Multimedia*. 533–542.
- [40] Fatemehsadat Mireshghallah, Vaishnavi Shrivastava, Milad Shokouhi, Taylor Berg-Kirkpatrick, Robert Sim, and Dimitrios Dimitriadis. 2021. Useridentifier: Implicit user representations for simple and effective personalized sentiment analysis. arXiv preprint arXiv:2110.00135 (2021).
- [41] Sounak Mondal, Zhibo Yang, Seoyoung Ahn, Dimitris Samaras, Gregory Zelinsky, and Minh Hoai. 2023. Gazeformer: Scalable, Effective and Fast Prediction of Goal-Directed Human Attention. arXiv preprint arXiv:2303.15274 (2023).
- [42] Yuya Moroto, Keisuke Maeda, Takahiro Ogawa, and Miki Haseyama. 2020. Few-shot personalized saliency prediction based on adaptive image selection considering object and visual attention. Sensors 20, 8 (2020), 2170.
- [43] Junting Pan, Elisa Sayrol, Xavier Giro-i Nieto, Kevin McGuinness, and Noel E O'Connor. 2016. Shallow and deep convolutional networks for saliency prediction. In Proceedings of the IEEE conference on computer vision and pattern recognition. 598–606.
- [44] Michael J Pazzani, Jack Muramatsu, Daniel Billsus, et al. 1996. Syskill & Webert: Identifying interesting web sites. In AAAI/IAAI, Vol. 1. 54-61.
- [45] Evan F Risko, Nicola C Anderson, Sophie Lanthier, and Alan Kingstone. 2012. Curious eyes: Individual differences in personality predict eye movement behavior in scene-viewing. Cognition 122, 1 (2012), 86–90.
- [46] Negar Sammaknejad, Hamidreza Pouretemad, Changiz Eslahchi, Alireza Salahirad, and Ashkan Alinejad. 2017. Gender classification based on eye movements: A processing effect during passive face viewing. Advances in cognitive psychology 13, 3 (2017), 232.
- [47] Hosnieh Sattar, Mario Fritz, and Andreas Bulling. 2020. Deep gaze pooling: Inferring and visually decoding search intents from human gaze fixations. Neurocomputing 387 (2020), 369–382.
- [48] Jude Shavlik, Susan Calcari, Tina Eliassi-Rad, and Jack Solock. 1998. An instructable, adaptive interface for discovering and monitoring information on the world-wide web. In Proceedings of the 4th international conference on Intelligent user interfaces. 157–160.
- [49] Ana Filipa Silva, Francisco Tomás González Fernández, et al. 2022. Differences in visual search behavior between expert and novice team sports athletes: A systematic review with meta-analysis. (2022).
- [50] Ekta Sood, Fabian Kögel, Florian Strohm, Prajit Dhar, and Andreas Bulling. 2021. VQA-MHUG: A gaze dataset to study multimodal neural attention in visual question answering. arXiv preprint arXiv:2109.13116 (2021).
- [51] Ekta Sood, Fabian Kögel, Philipp Müller, Dominike Thomas, Mihai Bâce, and Andreas Bulling. 2023. Multimodal Integration of Human-Like Attention in Visual Question Answering. In Proc. Workshop on Gaze Estimation and Prediction in the Wild (GAZE), CVPRW. 2647–2657. https://openaccess.thecvf.com/content/CVPR2023W/GAZE/papers/Sood_Multimodal_Integration_of_Human-Like_Attention_in_Visual_Question Answering CVPRW 2023 paper.pdf
- [52] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. The journal of machine learning research 15, 1 (2014), 1929–1958.
- [53] Florian Strohm, Ekta Sood, Sven Mayer, Philipp Müller, Mihai Bâce, and Andreas Bulling. 2021. Neural Photofit: gaze-based mental image reconstruction. In Proceedings of the IEEE/CVF International Conference on Computer Vision. 245–254.
- [54] Florian Strohm, Ekta Sood, Dominike Thomas, Mihai Bâce, and Andreas Bulling. 2023. Facial Composite Generation with Iterative Human Feedback. In Annual Conference on Neural Information Processing Systems. PMLR, 165–183.
- [55] Roel Vertegaal. 2002. Designing attentive interfaces. In Proceedings of the 2002 symposium on Eye tracking research & applications. 23-30.
- [56] Lennart Wachowiak, Peter Tisnikar, Gerard Canal, Andrew Coles, Matteo Leonetti, and Oya Celiktutan. 2022. Analysing eye gaze patterns during confusion and errors in human-agent collaborations. In 2022 31st IEEE International Conference on Robot and Human Interactive Communication (RO-MAN). IEEE, 224–229.
- [57] Dirk Walther and Christof Koch. 2006. Modeling attention to salient proto-objects. Neural networks 19, 9 (2006), 1395–1407.
- [58] Xiaodong Wu, Weizhe Lin, Zhilin Wang, and Elena Rastorgueva. 2020. Author2vec: A framework for generating user embedding. arXiv preprint arXiv:2003.11627 (2020).

- [59] Yanyu Xu, Shenghua Gao, Junru Wu, Nianyi Li, and Jingyi Yu. 2018. Personalized saliency and its prediction. *IEEE transactions on pattern analysis and machine intelligence* 41, 12 (2018), 2975–2989.
- [60] Yanyu Xu, Nianyi Li, Junru Wu, Jingyi Yu, Shenghua Gao, et al. 2017. Beyond Universal Saliency: Personalized Saliency Prediction with Multi-task CNN.. In IJCAI. 3887–3893.
- [61] Zhibo Yang, Lihan Huang, Yupei Chen, Zijun Wei, Seoyoung Ahn, Gregory Zelinsky, Dimitris Samaras, and Minh Hoai. 2020. Predicting goal-directed human attention using inverse reinforcement learning. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 193–202.
- [62] Alfred L Yarbus and Alfred L Yarbus. 1967. Eye movements during perception of complex objects. Eye movements and vision (1967), 171-211.
- [63] Bingqing Yu. 2018. Personalization of saliency estimation. McGill University (Canada).
- [64] Qian Zhao, Shuo Chang, F Maxwell Harper, and Joseph A Konstan. 2016. Gaze prediction for recommender systems. In Proceedings of the 10th ACM Conference on Recommender Systems. 131–138.
- [65] Wanjun Zhong, Duyu Tang, Jiahai Wang, Jian Yin, and Nan Duan. 2021. UserAdapter: Few-shot user learning in sentiment analysis. In Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021. 1484–1488.