

# Attention Flow: End-to-End Joint Attention Estimation

Ömer Sümer<sup>1</sup>, Peter Gerjets<sup>2</sup>, Ulrich Trautwein<sup>1</sup>, Enkelejda Kasneci<sup>1</sup>

<sup>1</sup> University of Tübingen <sup>2</sup> Leibniz-Institut für Wissensmedien  
Tübingen, Germany

{oemer.suemer, ulrich.trautwein, enkelejda.kasneci}@uni-tuebingen.de

p.gerjets@iwm-tuebingen.de

## Abstract

*This paper addresses the problem of understanding joint attention in third-person social scene videos. Joint attention is the shared gaze behaviour of two or more individuals on an object or an area of interest and has a wide range of applications such as human-computer interaction, educational assessment, treatment of patients with attention disorders, and many more. Our method, Attention Flow, learns joint attention in an end-to-end fashion by using saliency-augmented attention maps and two novel convolutional attention mechanisms that determine to select relevant features and improve joint attention localization. We compare the effect of saliency maps and attention mechanisms and report quantitative and qualitative results on the detection and localization of joint attention in the VideoCoAtt dataset, which contains complex social scenes.*

## 1. Introduction

Humans spend most of their lives interacting with each other. In public or private spaces such as squares, concert halls, cafes, schools, we share various aspects of everyday life with one another. Through new technologies and growing distractive effects of social media, we divide our attention and memory into separate themes and may have difficulties to focus our attention onto our primary task. In that regard, from both psychological and computer vision perspectives, understanding a person’s attentional focus and particular localization of joint attention present valuable opportunities.

Joint attention is very helpful in many different contexts. For example, in classroom-based learning, teachers who engage all students equally can enhance student achievement [12, 20, 30, 39]. To investigate this, educational researchers manually analyze student behaviours and especially the visual attention of students from video recordings of instructions and try to explain relationships between students’ and teachers’ behaviour in a very time-consuming way. Another

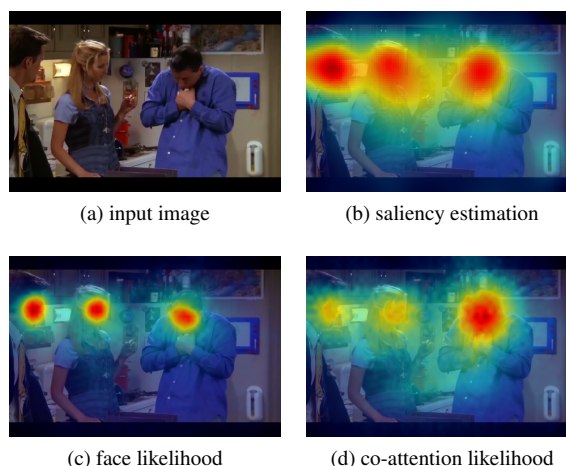


Figure 1. Sample of a social scene in (a), and the estimated saliency map using [21] in (b). Our method, *Attention Flow* takes only the input image in (a) and estimate the face likelihood (c) and the co-attention likelihood (d).

example is in the context of attention disorders or autism research. For instance, it has been shown that joint attention and engagement, particularly in early ages, can be taught using behavioural and developmental interventions [25]. Thus, computer vision-based, automated joint attention analysis can be instrumental in behavioural psychology to develop efficient training curricula for the treatment of children with disabilities. Another useful application is in the area of human-computer interaction and especially interaction with autonomous systems. For example, robots can infer gaze direction in case of a single person or joint attention in groups and turn their heads into that direction. Such information could be further used by robots to augment their collaboration with humans [2, 3].

Although an automated analysis of joint attention might be beneficial for a variety of applications, related work in the domain of computer vision is still quite limited. Few works addressed a similar problem, namely social saliency

in first and third-person view [27, 28, 29, 37]. Also, there are examples of joint attention in human-robot interaction [2, 9, 23, 24]. Whereas mapping gaze directions to a common plane [33] is a promising option in controlled settings, it does not work in more challenging multimedia data. A recent study [10] collected a large video dataset, which we use in this study, and proposed a spatiotemporal neural network to estimate shared attention. Even though we deal with the same problem, we prefer to use the term of joint attention, since shared attention, from a psychological perspective, includes further underlying cognitive processes and does not necessitate joint gaze.

In this work, we propose a new approach that relates saliency and joint attention to estimate locations of joint attention in third person images or videos. Simply explained, saliency is an estimation of fixation likelihood on an image. In fact, due to the limited capacity of our visual system, we, by the help of an attentional mechanism, focus on the most relevant parts of a scene that are more distinctive than the remaining. In essence, it is how our eye movements process a scene, by employing various eye movements (such as saccade and fixations) and visual search which is guided by various bottom-up and top-down processes. Eye tracking-based saliency information has supported many computer vision tasks such as object detection [26], zero-shot image classification [19], and image/video captioning [8, 43].

Figure 1 shows a sample of our approach. Despite the usefulness of saliency maps, they do not necessarily represent the visual focus of people in the scene. However, during the training time, we exploit saliency maps to encode contextual information and create pseudo attention maps by combining them with face locations and their joint attention point and learn to predict these likelihoods. Then, during the test time, we can summarise the attentional focus of people in given third-person social images or videos.

The main contributions of this paper are as follows:

1. It formulates the problem of inferring joint attention as end-to-end training. Thereby, *Attention Flow* works without additional dependencies such as face/head detection, region proposals, or saliency estimation.
2. It explicitly learns saliency and joint attention of a high-level inference task using saliency augmented pseudo attention maps and Attention Flow network with channel-wise and spatial attention mechanisms.
3. Experimental results verify the performance of our approach on large-scale social videos, namely the Video-CoAtt dataset [10]. We also present a comparative ablation analysis of saliency and attention modules.

## 2. Related Work

First, we review related research on gaze following and joint attention. Then, we will discuss saliency estimation and attention modules as we utilized them in our approach

to infer the joint attention.

**Gaze Following:** Recasens *et al.* [31] proposed a neural network which predicts the locations being gazed at in a convolutional neural network using head location, an image patch from head location, and an entire image. They also created a large-scale dataset where persons’ eye and gaze locations were annotated. They later extended their work to use eye locations in a video frame and to predict gazed location in future frames [32]. Gorji *et al.* [13] used a similar approach to [31]; however, they leveraged gaze information to boost saliency estimation and did not report gaze following results.

Recently, Chong *et al.* [7] proposed a method to train gaze following, head pose and gaze tasks based on a multi-task learning approach by optimizing several losses on different tasks and datasets. They also included *outside* of the frame labels and predicted visual attention. Nevertheless, their approach estimates a single person’s visual attention, not joint attention.

**Joint vs. Shared Attention:** Joint attention is a social interaction that can occur in the forms of dyadic (looking at each other) or triadic ways (looking at each other and an object). Previous research shows that infants can discriminate between dyadic and joint attention interactions already by the age of 3 months [36]. Joint attention is crucial for language learning and imitative learning [1, 22]. In contrast to joint attention, shared attention does not require co-attending physically or by gaze. For instance, co-attending a television broadcast when looking at another point can be an example of shared attention. An observer can understand shared attention by using cues from the environment [35]. Shared attention is more related to the underlying cognitive processes, whereas joint attention is dyadic and triadic gaze oriented. Thus, in the following we will use the term of joint attention since computer vision relies on seen visual cues.

**Joint Attention:** Looking into studies on the analysis of attention in social interactions, [27] localized head-mounted cameras in 3D using structure from motion and triangulated joint attention. Later, they proposed a geometric model between joint attention and social formation captured from first and third person views [29]. These works are noteworthy; however, they depend on first-person views and thus cannot be applied in unconstrained third view images and videos. Also, they aim to predict only proximity of joint attention (social saliency) and cannot present a good understanding of joint attention.

**Saliency Estimation:** Saliency is a measure of spatial importance, and it characterizes the parts of the scene which stand out relative to other parts. Being salient can depend on low-level features such as luminance, color, texture, high-level features such as objectness, task-driven factor, and center bias phenomenon. In the literature of saliency estimation, two approaches exist: (a) bottom-up methods,

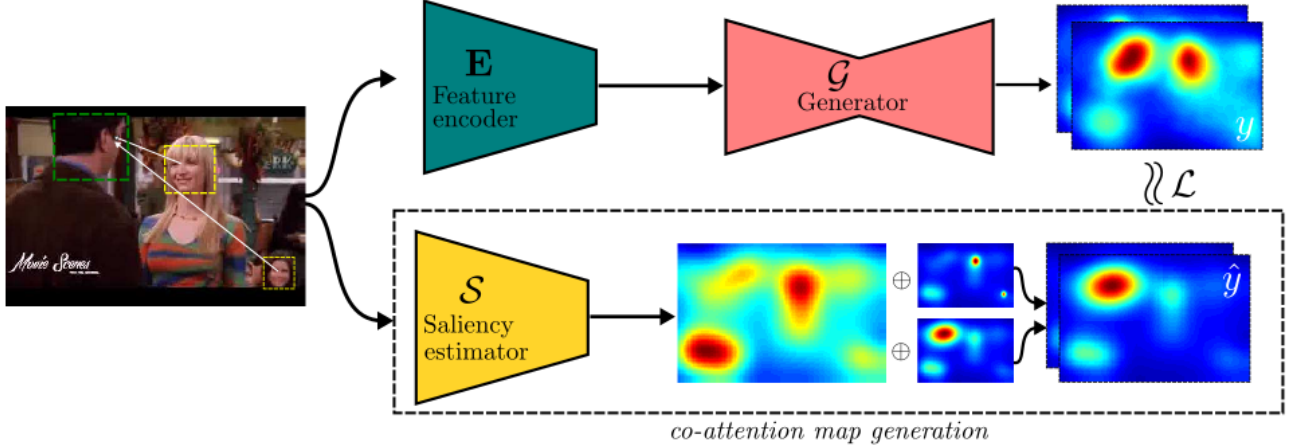


Figure 2. *Overview of our Attention Flow* Our method is composed of three modules, (i) feature encoder, (ii) attention flow generator, and (iii) saliency-based ground truth generation. It estimates a two-channel heatmap, which encodes faces and their co-attention likelihood in the scene.

which aim to combine relevant information without prior knowledge of the scene, and (b) top-down methods which are more goal-oriented [5]. Availability of large-scale attention datasets and deep learning approaches have surpassed all previous psychological and computational methods. Based on these recent studies, we know that humans look at humans, faces, objects, texts [6] and also emotional content [11]. The joint attention of humans in the scene is also noticeable. For this reason, we will leverage saliency information to learn joint attention.

**Attention Mechanism:** Computer-based estimation of attention can also be approached by means of machine-learning techniques, where models, with the help of spatial or temporal attention mechanisms, are able to learn where, when, and what to attend. The use of First use cases are machine translation [4], image captioning [41], and action classification [34].

Looking into attention mechanisms in images, Wang *et al.* [38] incorporated attention modules into an encoder-decoder network and performed well in an image classification task. Their method learns attention jointly in 3D. Another recent work exploited inter-channel relationships. In Squeeze-and-Excitation blocks, they utilized global average-pooled features to perform a channel-wise calibration [17]. Recently, Woo *et al.* [40] proposed a convolutional attention module that leverages channel and spatial relations separately.

The common point of these works is that they address classification tasks by the use of spatial, temporal, or channel-wise attention. In contrast, we propose novel convolutional attention mechanisms for two purposes: the first is to learn feature selection along the channel dimension of a learned representation, and secondly, to guide a regression network to focus on more relevant areas in the spatial dimension. Instead of an architectural block in a classifica-

tion task as in [17], we utilize these blocks to benefit from learned features better by applying an adaptive feature selection and apply a further refinement on top of the heatmap generation module.

### 3. Method

Our approach aims to infer joint attention in third person social videos, where two or more people look at another person or object. Figure 2 shows an overview of our workflow.

For a given social image or video frame, we estimate a two-channel likelihood distribution, called *Attention Flow*. One channel represents faces in the scene, whereas the second channel is the likelihood of joint attention. In our workflow, raw images can be considered as a fusion of social presence in the scene and the center of joint attention. Figure 2 depicts an example prediction of our approach. Our Attention Flow network takes only raw images and detects faces and their respective co-attention locations without depending on any other information. In this section, we will describe (1) the creation of pseudo-attention maps (§3.1), which are augmented by saliency estimation; (2) learning and inference (§3.2) by our Attention Flow network using attention mechanisms, and provide (3) implementation details (§3.3).

#### 3.1. Saliency Augmented Pseudo-Attention Maps

Consider persons interacting with each other in a social scene. The question we address is how to infer their visual attention focus from a third person’s view? Probably the most accurate way to obtain this information would be by employing mobile eye trackers or through gaze estimation based on several high-resolution field cameras. For the majority of use cases in our daily lives, where such

equipment cannot be employed, it would be very useful to be able to retrieve such information solely based on images or video material. For this reason, we first compute pseudo-attention maps by leveraging saliency estimation.

More specifically, for an input image  $I$ , we have a number of detected head locations  $(x_i, y_i, w_i, h_i)$ , where  $i = 1, 2, \dots, n$  and  $n \geq 0$ . To model social presence and the respective co-attention location we use Gaussian distributions. For a head detection or co-attention bounding box, this distribution is defined as

$$\mathcal{G}(x + \delta x, y + \delta y) = \begin{cases} \exp\{-\frac{x^2+y^2}{2\sigma^2}\} & \|\delta x\| \leq w, \|\delta y\| \leq h \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

Then, we combine head locations and co-attention maps with the estimated saliency maps, which is a precursor of observer's attention. Augmented by estimated saliency maps  $\mathcal{S}$ , the created pseudo-attention maps can be formalized as follows:

$$\begin{aligned} \mathcal{H}_1 &= \alpha \log\left(\sum_{i=1, \dots, n} \mathcal{G}_{f_i}\right) + \beta \log(\mathcal{S}) \\ \mathcal{H}_2 &= \alpha \log(\mathcal{G}_{coatt}) + \beta \log(\mathcal{S}) \end{aligned} \quad (2)$$

In this way, we suppress the saliency to lower values and ensure that if there are detected faces in the scene, they and their respective co-attention point will correspond to the maximum values of pseudo-attention maps in the first and second channels.

By employing saliency estimation in our method, we leverage the information of relative importance of the regions which can be also salient for the persons in the scene. Thus, it prevents unreliable training samples, where the same object can appear as a co-attention point or zero when we use only  $\mathcal{G}_f$  and  $\mathcal{G}_{coatt}$ .

### 3.2. Attention Flow Network

Our model aims to solve three problems simultaneously: (1) to locate faces in the given image, (2) to detect whether joint attention exists or not, and (3) to predict the location of joint attention.

As input we only use the raw images instead of any other computational blocks, such as face detector, object detector or proposal networks. In this way, our Attention Flow network can be used to retrieve images or videos according to their social context in an efficient and fast way. The two-channel saliency augmented pseudo-attention maps are a compressed form of these objectives and provide all necessary information. In images which do not contain faces,

the first channel of the attention map will give a lower likelihood, and they can be easily omitted from the further attention analysis.

In case of two or more persons in the scene, the first channel will represent the locations of their faces, whereas the second channel will be either estimated saliency or joint attention. Since pseudo-attention maps are a weighted summation of saliency estimation and joint attention, the typical values of maximum points are informative about the presence of joint attention. Therefore, learning pseudo-attention maps enables both detection and localization tasks simultaneously.

As it can be seen in Figure 2, we first extract a visual representation of the scene using a pre-trained encoder network on object classification tasks. Since inferring joint attention is a complex problem even for humans, we leverage from an encoder to understand the visual focus of the persons in the image and for better generalization. The following block is a generator network, which learns attention maps from encoded representations. In order to avoid undesired outcomes of rescaling, we preserve the original aspect ratio in the input image and prefer fully convolutional architectures in both encoder and generator networks.

As a loss function, we use the Mean Squared Error (MSE) between the predicted attentions maps  $\mathcal{H}$  and ground truth pseudo saliency maps  $\mathcal{H}$  (created as described in §3.1):

$$\mathcal{L}_{MSE} = \frac{1}{H \cdot W} \|G(E(I)), \mathcal{H}\|^2 \quad (3)$$

When compared to other vision tasks such as object detection, segmentation or categorization, localizing the joint attention is a very complex task because the same region, i.e., a face or an object, can be the co-attention point in a scene, but shortly for a short period of time, it might not be true. In order to deal with these situations, our Attention Flow network can be guided towards the more relevant regions. For this purpose, we propose two novel attention mechanisms, namely *channel-wise* and *spatial*, and investigate their efficiency in the localization of joint attention. Figure 3 shows these attention mechanisms.

The encoder output is typically the output of a convolutional network which preserves spatial information in a reduced resolution and contains a higher dimension in the channel. The combination of these feature maps decides whether objects are present in the image. Using the complete encoded representation is redundant. According to the context, some channels can have more importance in the representation of the scene. Channel-wise convolutional attention performs a feature selection by weighting channels according to their contribution to the task.

On the other hand, spatial attention works as a refinement on top of the final joint attention estimations. In contrast to the spatial attention mechanisms in classification,

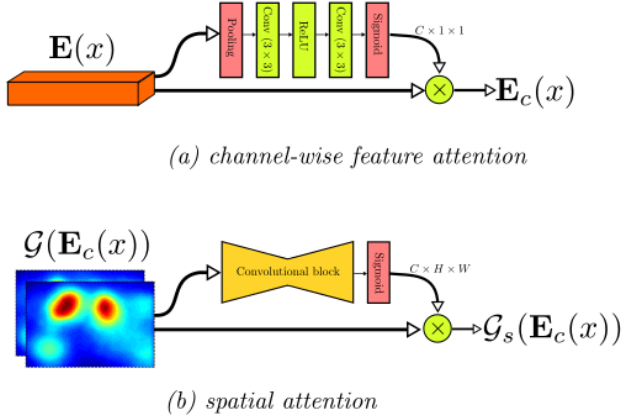


Figure 3. Channel-wise feature attention and convolutional spatial attention blocks.

which works as an importance map to maximize class activations, our spatial attention augments a heatmap regression task.

### 3.3. Implementation Details

The Attention Flow network is composed of three main modules: encoder, generator, and co-attention map generation blocks. In order to exploit the knowledge of large-scale object classification tasks, we use a pre-trained ResNet-50 [15] as an encoder. Our final estimation is an attention map and needs to preserve spatial relations as much as possible. Thus, we prefer dilated residual architecture, DRN-A-50 [42] trained on ImageNet and keep the resolution 1/8 at the output of the encoder.

As a generator, we used 9 residual blocks with instance normalization. It takes inputs in the number of feature channels (2048) and outputs 2-channel attention maps. Then, linear upsampling (x8) is applied.

The last block is co-attention map generation, and it is used only in training time to produce ground truth attention maps as described in §3.1. To estimate saliency, we used Deep Gaze II [21]. Similar to other data-driven saliency estimation methods, Deep Gaze II makes use of different level of features and has an understanding of objectness. It helps us to reduce the number of potential locations where joint attention might exist.

The layers of channel-wise feature attention are depicted in Figure 3(a). On the other hand, in the convolutional block of spatial attention, we used a small residual network that contains 3 residual blocks. As we applied spatial attention at 1/8 resolution before upsampling, it does not introduce an extensive computational cost to the entire workflow.

At training time, we used a SGD solver with a learning rate of 0.01 in the generator block. In feature encoder, we either lock the pre-trained parameters or applied *fine-tuning* by a 10 times reduced learning rate.

## 4. Experiments

In this section, we first define the used dataset and performance metrics. Then, we report the ablation studies on the use of saliency estimation to create attention maps and the effect of attention mechanisms and evaluate our approach on the VideoCoAtt dataset in comparison to related approaches.

### 4.1. Experimental Setup

To evaluate our approach on joint attention estimation, we used the Video Co-Attention dataset [10], which is currently, to the best of our knowledge, the only available dataset on a joint attention task. The dataset contains 380 RGB video sequences from 20 different TV shows in the resolution of  $320 \times 480$ . There are 250,030 frames in the training set, 128,260 frames in the validation set, and 113,810 frames in the testing set. Each split comes from different TV shows, and the dataset includes varying human appearances and formation.

There are two tasks: detection and localization of joint attention. Some images might not contain human bodies or faces. In images with social content, subjects' attentional focus can be different. In the detection task, we report overall prediction accuracy in the test set of VideoCoAtt. On the other hand, localization is evaluated on the test images with joint attention locations.  $L_2$  distance in the input resolution will be used.

By adopting the evaluation procedure from [10], we use the Structured Edge Detection Toolbox [44] to generate bounding box proposals. In the location, where our method predicts joint attention, we apply a Non-Maximum Suppression (NMS) and take the one that intersects greatest with our predicted estimation. It should be noted that our approach can locate the center of joint attention. Thus, in order to make a fair comparison with state-of-the-art methods, and we used the bounding box proposal.

Furthermore, there may be no joint attention or more than one joint attention location in an image. In order to learn the detection and localization of joint attention at the same time, we learn by all types of images without social context (body or faces), with social context but without any joint attention, and one or more joint attention. According to Eq. 2, we limit the values of saliency to some range by natural logarithm and a scale factor. Thus, our trained network's prediction can be joint attention if and only if the predicted likelihood is greater than a threshold.

### 4.2. Results and Analysis

**Saliency and joint attention** Saliency models the attention of a third person who observes a video or image. On the other hand, joint attention analysis aims to understand from the perspectives of persons in these visual content. Due to





Figure 4. Example daily life scenes from VideoCoAtt dataset [10] and their respective saliency estimations using Deep Gaze II [21]. The focus of shared visual attention does not necessarily need to be the most salient region, but contains auxiliary information to localize joint attention.

the geometric difference between the viewpoints and human behavior in social scenes, the most salient part of images may not be the focus of persons’ attention. Thus, we investigate how the co-attention locations are salient for different saliency estimation methods.

We tested four saliency estimation methods, Itti and Koch [18], GBVS [14], Signature [16], and Deep Gaze 2 [21]. The first three were chosen as representatives of classical computational saliency methods, whereas Deep Gaze 2 is a data-driven approach that depends on pre-trained feature representations on image classification. Deep Gaze 2’s mean saliency value in co-attention bounding boxes of the training images, 96% of the time, are above the mean saliency value of images, whereas it is the case in 44%, 71% and 77% for Itti & Koch, GBVS, and Signature.

In most cases, persons in the scene interact with either another person or an object. We regard that a data-driven Deep Gaze 2 can result in higher saliency in co-attention regions as it leverages a representation trained on object classification. Thus, we prefer Deep Gaze 2 when creating pseudo attention maps (§3.1).

Figure 4 shows some sample images and their estimated saliency maps using Deep Gaze 2 [21], respectively. These samples show us that the most salient regions do not necessarily contain the possible joint attention in the social images. However, they are a precursor of observer’s attention who gaze at images.

A tiny visual change in the image can cause a big change in the presence and location of joint attention. This is the main reason why we leverage the saliency information. The “raw image” results in [10] also validate our assumption. One can suppose the use of saliency as introducing noise, however, starting from the attention of observer and guiding the attention of the network towards understanding the

attention of people inside the scene is a reasonable solution and also makes the problem learnable.

**Use of attention mechanisms** Our Attention Flow network learns joint attention by using a pre-trained representation and a generator as a regression task with mean square loss. To supplement it, we proposed two attention mechanisms. In contrast to existing attention mechanisms in the literature, such as temporal in videos or text data, or spatial in image categorization, we use two novel convolutional attention blocks for feature selection and regression tasks. We evaluate their performance on the joint attention localization task.

The output of the dilated residual network that we used as an encoder is 1/8 resolution of the input and its channel size is 2048. The channel-wise attention module, first applies  $(4 \times 6)$  average pooling, two convolutional layers whose kernel sizes are  $3 \times 3$ ,  $3 \times 2$  with a stride of 2 and 1, respectively. Their channel sizes are 512 and 2048, respectively. The final output is in the size of  $C \times 1 \times 1$  and the original encoder output is channel-wise multiplied by these importance weights.

Table 1 shows the results of joint attention localization over the test set of VideoCoAtt dataset. We first tested the following options: To use the encoder as pre-trained features (no learning), to train the encoder in the same learning rate as the generator, and finetuning the encoder by a reduced ( $\times 0.1$ ) learning rate. Channel-wise attention aims to apply feature selection in a learn representation. Thus, we freeze the encoder when training channel-wise attention and generator jointly. This approach reduces the mean  $L_2$  distance by 10.92 and 6.88 pixels in comparison with no learning and finetuning, respectively.

Looking into our spatial attention, we applied spatial attention to the output of the generator in 1/8 resolution ( $40 \times 60$ ) before linear upsampling. Spatial attention module takes the estimation of joint attention maps ( $2 \times H \times W$ ) and learns a spatial importance on top to better localize the co-attention point. In spatial attention, we use a  $3 \times 3$  convolutional layer (64), batch normalization, a residual bottle-

Method	$L_2$ distance
Attention Flow	
$E_{(lr=0)}$	73.77
$E_{(base\_lr)}$	70.47
$E_{(finetune)}$	69.72
<i>channel – wise attention</i>	62.84
<i>spatial attention</i>	65.70

Table 1. The effect of attention mechanisms in localization of joint attention over the test set of VideoCoAtt dataset.



Figure 5. **Qualitative results of Attention Flow** Bounding boxes on sample test images (green) show the ground truth attentional focus. The second and third columns are our estimated Attention Flow. In the third column, we also depicted the estimated bounding boxes (cyan). Figure best viewed in color.

neck module and final convolutional layer to reduce channel size back to 2. Then, a sigmoid activation is applied and the previous predictions are weighted.

Before using attention mechanisms, we show how accurate we can localize co-attention bounding boxes based on our baseline approach that is depicted in Figure 2. After creating saliency guided pseudo-attention maps that we use as the label, our Attention Flow network has two trainable blocks: an encoder, and a generator. The encoder is initialized by ImageNet trained weights. Then, we compared the following three cases: freeze the encoder and train only generator ( $E_{(lr=0)}$ ), train encoder and generator jointly ( $E_{(base\_lr)}$ ), and learn encoder by transfer learning with a reduced learning rate and train generator from scratch ( $E_{(finetune)}$ ). As it can be seen in Table 1, transfer learning performs better than the approaches mentioned above and achieves an  $L_2$  distance of 69.72.

Channel-wise attention, which is used between encoder and generator, can predict joint attention with a mean distance of 62.84, whereas spatial attention after the generator gives 65.70. Both attention mechanisms improve joint attention localization by 4.02 and 6.88 points with respect to our baseline network with encoder fine-tuned in Table 1. The better performance of our channel-wise attention approach indicates that feature selection on top of deep learning features plays an important role. Weighting features per channel improves their potential as a scene descriptor.

Table 2 shows our results in comparison with other methods in detection and localization of joint attention. **Ran-**

**dom** is acquired by drawing a Gaussian heatmap with a random mean and variance. **Fixed Bias** uses joint attention bias in the TV shows (averaged over the VideoCoAtt dataset) to sample predictions. An alternative to joint attention is to make prediction per person using **Gaze Follow** [31] and combine their attention likelihoods. Other methods are from the reference of VideoCoAtt dataset [10] and grouped into two categories: single frame and temporal models. All of these methods [10] depend on head detection bounding boxes, region proposal model or saliency estimation even in test time. In terms of used modalities, our approach is similar to their “**Raw Image**” approach.

Our Attention Flow network with a channel-wise attention detects joint attention with an accuracy of 78.1% over

Model	Prediction Acc.	$L_2$ distance
Random	50.8%	286
Fixed Bias	52.4%	122
Gaze Follow [31]	58.7%	102
Raw Image [10]	52.3%	188
Only Gaze [10]	64.0%	108
Gaze+RP [10]	68.5%	74
Gaze+Saliency+LSTM [10]	66.2%	71
Gaze+RP+LSTM [10]	71.4%	<b>62</b>
<b>Ours (w channel-wise att.)</b>	<b>78.1%</b>	62.84

Table 2. **Quantitative evaluation results** with Prediction Accuracy and  $L_2$  Distance over the test set of VideoCoAtt dataset.

the entire test set of VideoCoAtt. Furthermore, it localizes co-attention bounding boxes with  $L_2$  distance of 62.84. Our method performs significantly better than [10]’s single frame with region proposals and gaze estimation. Furthermore, our approach is on par with **Gaze+RP+LSTM** and outperforms it in terms of prediction accuracy by 6.7%.

We should note that our model makes this improvement without using any head pose/gaze estimation branch, region proposal maps, and also temporal information. [9]’s models with LSTM leverages 20-30-frame length sequences to improve and smooth prediction performance. We focused on learning an end-to-end model by using only single raw frames. Therefore, as in Table 2, our model’s performance (78.1% and 62.84) is far beyond [31] and [10]’s single frame approaches which perform at best, **Gaze+RP**, 68.5% and 74 in joint attention detection and localization, respectively.

Figure 5 depicts the qualitative results of our Attention Flow network on several test images of VideoCoAtt. The ground truth co-attention locations are shown in green rectangles. Estimated face and co-attention likelihoods are overlaid on images. The first channel of our attention maps successfully locates both frontal and side faces. Looking into co-attention estimation, predictions in groups with 3-4 persons are very good. Even though their distances from ground truths are not very large, the last two examples (on the right) are relatively broad. This is due to the difficulty of scenes and a wider angle of view.

Another point that we should address is the distribution of social formations in the VideoCoAtt dataset. As the dataset is composed of acted scenes mostly from the TV shows, it does not represent the real-world formations such as in learning situations or group work. In addition, when we inspect the failures, we observed that most of them were from complicated cases where many people interest each other. Their faces were far from the camera and difficult to determine their activities (i.e., last two samples on the right side of Figure 5). The possible direction in joint attention analysis can be to create training corpus specialized in the desired applications such as group work analysis, therapeutic situations, or children’s gaze behaviors.

## 5. Conclusion

This study addressed a recently proposed problem, inferring joint attention in third person social videos. Our Attention Flow network infers joint attention based on only raw input images. Without using any temporal information and other dependencies such as a face detector or head pose/gaze estimator, we detect and localize joint attention better than the previous approaches. We create pseudo-attention maps by leveraging saliency information to better detect and localize joint attention. Furthermore, we propose two new convolutional attention blocks for feature selection

and attention map localization. As inferring joint attention in an end-to-end fashion necessitates a high-level inference, increasing the amount of training data or the network depth will not help. We should note that these attention mechanisms, particularly channel-wise attention blocks for feature selection, are highly essential to select useful features from learned representations and improve localization performance of a heatmap regressor.

Understanding of joint attention by use of computer vision can help in a wide range of applications such as educational assessment, human-computer interactions, and therapy for attention disorders and as a future work we extend our approach to specialize in these applications and use as a tool for human behavior understanding.

**Acknowledgements** Ömer Sümer is a doctoral student at the LEAD Graduate School & Research Network [GSC1028], funded by the Excellence Initiative of the German federal and state governments. This work is also supported by Leibniz-WissenschaftsCampus Tübingen “Cognitive Interfaces”.

## References

- [1] D. A. Baldwin. Early referential understanding: Infants’ ability to recognize referential acts for what they are. *Developmental Psychology*, 29:832–843, 09 1993.
- [2] S. Andrist, B. Mutlu, and A. Tapus. Look like me: Matching robot personality via gaze to increase motivation. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems, CHI ’15*, pages 3603–3612, New York, NY, USA, 2015. ACM.
- [3] R. M. Aronson, T. Santini, T. C. Kübler, E. Kasneci, S. Srinivasa, and H. Admoni. Eye-hand behavior in human-robot shared manipulation. In *Proceedings of the 2018 ACM/IEEE International Conference on Human-Robot Interaction, HRI ’18*, page 4–13, New York, NY, USA, 2018. Association for Computing Machinery.
- [4] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. *arXiv e-prints*, abs/1409.0473, Sept. 2014.
- [5] A. Borji and L. Itti. State-of-the-art in visual attention modeling. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(1):185–207, Jan 2013.
- [6] Z. Bylinskii, A. Recasens, A. Borji, A. Oliva, A. Torralba, and F. Durand. Where should saliency models look next? In *Computer Vision – ECCV 2016*, pages 809–824, Cham, 2016. Springer International Publishing.
- [7] E. Chong, N. Ruiz, Y. Wang, Y. Zhang, A. Rozga, and J. M. Rehg. Connecting gaze, scene, and attention: Generalized attention estimation via joint modeling of gaze and scene saliency. In *The European Conference on Computer Vision (ECCV)*, September 2018.
- [8] M. Cornia, L. Baraldi, G. Serra, and R. Cucchiara. Predicting human eye fixations via an lstm-based saliency attentive



- model. *IEEE Trans. Image Processing*, 27(10):5142–5154, 2018.
- [9] R. R. da Silva and R. A. F. Romero. Modelling shared attention through relational reinforcement learning. *Journal of Intelligent & Robotic Systems*, 66(1):167–182, Apr 2012.
- [10] L. Fan, Y. Chen, P. Wei, W. Wang, and S.-C. Zhu. Inferring shared attention in social scene videos. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [11] S. Fan, Z. Shen, M. Jiang, B. L. Koenig, J. Xu, M. S. Kankanhalli, and Q. Zhao. Emotional attention: A study of image sentiment and visual attention. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [12] P. Goldberg, Ö. Sümer, K. Stürmer, W. Wagner, R. Göllner, P. Gerjets, E. Kasneci, and U. Trautwein. Attentive or not?: Toward a machine learning approach to assessing students’ visible engagement in classroom instruction. *Educational Psychology Review*, 2019.
- [13] S. Gorji and J. J. Clark. Attentional push: A deep convolutional network for augmenting image saliency with shared attention modeling in social scenes. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3472–3481, July 2017.
- [14] J. Harel, C. Koch, and P. Perona. Graph-based visual saliency. In B. Schölkopf, J. C. Platt, and T. Hoffman, editors, *Advances in Neural Information Processing Systems 19*, pages 545–552. MIT Press, 2007.
- [15] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 770–778, 2016.
- [16] X. Hou, J. Harel, and C. Koch. Image signature: Highlighting sparse salient regions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(1):194–201, Jan 2012.
- [17] J. Hu, L. Shen, and G. Sun. Squeeze-and-excitation networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [18] L. Itti, C. Koch, and E. Niebur. A model of saliency-based visual attention for rapid scene analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(11):1254–1259, Nov 1998.
- [19] N. Kaessli, Z. Akata, B. Schiele, and A. Bulling. Gaze embeddings for zero-shot image classification. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [20] A. Kent. Synchronization as a classroom dynamic: A practitioner’s perspective. *Mind, Brain, and Education*, 7(1):13–18, 2013.
- [21] M. Kümmerer, T. S. A. Wallis, and M. Bethge. Deepgaze II: reading fixations from deep features trained on object recognition. *CoRR*, abs/1610.01563, 2016.
- [22] D. S. Murray, N. A. Creaghead, P. Manning-Courtney, P. K. Shear, J. Bean, and J.-A. Prendeville. The relationship between joint attention and language in children with autism spectrum disorders. *Focus on Autism and Other Developmental Disabilities*, 23(1):5–14, 2008.
- [23] Y. Nagai, M. Asada, and K. Hosoda. Learning for joint attention helped by functional development. *Advanced Robotics*, 20(10):1165–1181, 2006.
- [24] Y. Nagai, K. Hosoda, A. Morita, and M. Asada. A constructive model for the development of joint attention. *Connection Science*, 15(4):211–229, 2003.
- [25] L. B. Olswang, P. Dowden, J. Feuerstein, K. Greenslade, G. L. Pinder, and K. Fleming. Triadic gaze intervention for young children with physical disabilities. *Journal of Speech, Language, and Hearing Research*, 57(5):1740–1753, 2014.
- [26] D. P. Papadopoulos, A. D. F. Clarke, F. Keller, and V. Ferrari. Training object class detectors from eye tracking data. In D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, editors, *Computer Vision – ECCV 2014*, pages 361–376, Cham, 2014. Springer International Publishing.
- [27] H. S. Park, E. Jain, and Y. Sheikh. 3d social saliency from head-mounted cameras. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 422–430. Curran Associates, Inc., 2012.
- [28] H. S. Park, E. Jain, and Y. Sheikh. Predicting primary gaze behavior using social saliency fields. In *2013 IEEE International Conference on Computer Vision*, pages 3503–3510, Dec 2013.
- [29] H. S. Park and J. Shi. Social saliency prediction. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4777–4785, June 2015.
- [30] L. Prieto, K. Sharma, L. Kidzinski, M. Rodriguez-Triana, and P. Dillenbourg. Multimodal teaching analytics: Automated extraction of orchestration graphs from wearable sensor data. *Journal of Computer Assisted Learning*, 34(2):193–203, 2018.
- [31] A. Recasens, A. Khosla, C. Vondrick, and A. Torralba. Where are they looking? In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 199–207. Curran Associates, Inc., 2015.
- [32] A. Recasens, C. Vondrick, A. Khosla, and A. Torralba. Following gaze in video. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 1444–1452, Oct 2017.
- [33] T. Santini, T. Kübler, L. Draghetti, P. Gerjets, W. Wagner, U. Trautwein, and E. Kasneci. Automatic mapping of remote crowd gaze to stimuli in the classroom. In *Eye Tracking Enhanced Learning (ETEL2017)*, 09 2017.
- [34] S. Sharma, R. Kiros, and R. Salakhutdinov. Action recognition using visual attention. *CoRR*, abs/1511.04119, 2015.
- [35] G. Shteynberg. Shared attention. *Perspectives on Psychological Science*, 10(5):579–590, 2015. PMID: 26385997.
- [36] T. Striano, V. M. Reid, and S. Hoehl. Neural mechanisms of joint attention in infancy. *European Journal of Neuroscience*, 23(10):2819–2823, 2006.
- [37] O. Sümer, P. Goldberg, K. Stürmer, T. Seidel, P. Gerjets, U. Trautwein, and E. Kasneci. Teachers’ perception in the classroom. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2018.

- [38] F. Wang, M. Jiang, C. Qian, S. Yang, C. Li, H. Zhang, X. Wang, and X. Tang. Residual attention network for image classification. *CoRR*, abs/1704.06904, 2017.
- [39] K. Watanabe. Teaching as a dynamic phenomenon with interpersonal interactions. *Mind, Brain, and Education*, 7(2):91–100, 2013.
- [40] S. Woo, J. Park, J.-Y. Lee, and I. So Kweon. Cbam: Convolutional block attention module. In *The European Conference on Computer Vision (ECCV)*, September 2018.
- [41] 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 F. Bach and D. Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 2048–2057, Lille, France, 07–09 Jul 2015. PMLR.
- [42] F. Yu, V. Koltun, and T. Funkhouser. Dilated residual networks. In *Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [43] Y. Yu, J. Choi, Y. Kim, K. Yoo, S.-H. Lee, and G. Kim. Supervising neural attention models for video captioning by human gaze data. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [44] C. L. Zitnick and P. Dollár. Edge boxes: Locating object proposals from edges. In D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, editors, *Computer Vision – ECCV 2014*, pages 391–405, Cham, 2014. Springer International Publishing.