

Beyond F-formations: Determining Social Involvement in Free Standing Conversing Groups from Static Images

Lu Zhang^{1,2} and Hayley Hung¹

¹Delft University of Technology, Mekelweg 2, Delft, Netherlands,
{lu.zhang, h.hung}@tudelft.nl

²University of Twente, Drienerlolaan 5, Enschede, Netherlands

Abstract

In this paper, we present the first attempt to analyse differing levels of social involvement in free standing conversing groups (or the so-called F-formations) from static images. In addition, we enrich state-of-the-art F-formation modelling by learning a frustum of attention that accounts for the spatial context. That is, F-formation configurations vary with respect to the arrangement of furniture and the non-uniform crowdedness in the space during mingling scenarios. The majority of prior works have considered the labelling of conversing group as an objective task, requiring only a single annotator. However, we show that by embracing the subjectivity of social involvement, we not only generate a richer model of the social interactions in a scene but also significantly improve F-formation detection. We carry out extensive experimental validation of our proposed approach by collecting a novel set of multi-annotator labels of involvement on the publicly available Idiap Poster Data; the only multi-annotator labelled database of free standing conversing groups that is currently available.

1. Introduction

In recent years, the analysis of mingling scenarios has received growing attention. The potential of studying social patterns of behaviour in visual scenes has great potential with the recent advances in social signal processing [21]. Potential applications include enabling robots to approach a group and offer assistance in a socially intelligent manner [18], or social surveillance [3], image interpretation or retrieval [14].

A major challenge in visual scene interpretation is addressing the problem of bridging the semantic gap [14], which defines the disconnect between information that can be extracted from the pixels in an image and how a human might interpret its contents. Traditionally, this gap has been

attributed to the mapping of imagery data to objective interpretations such as the labelling of objects or activities in a scene. However, in recent years, scene analysis has started to consider more complex and subjective concepts such as safety [11] or ambiance [12]. Similarly, in the area of social surveillance [3], researchers have been trying to ascribe social meaning to social scenes. However, unlike conventional scene analysis, social surveillance bridges a more complex semantic gap that associates observable behavioural cues to social phenomena. We call this the *social semantic gap*. Since social phenomena are extremely complex, it is desirable to use findings from social psychology to help inform how visually observed behaviours could be linked to social phenomena to help bridge the gap in an informed manner.

Given the great advances already in person tracking and orientation detection, we focus on how these solutions can be used as behavioural input for bridging the social semantic gap. Specifically, we approach the novel problem of detecting *associates* of conversing groups (or the so-called F-formations). F-formations are defined by social psychology theory as [8]; as a spatial organization of people gathered for conversation where each member has an equal ability to sense all other members. The so-called *associates* of F-formations are defined by psychologists as people who are attached to an F-formation but do not have the same status as full members (see Figure 1 (a)).

To the best of our knowledge, state-of-the-art methods for F-formation detection [6, 2, 13, 19, 20] have made three simplifying assumptions. First, each individual is assumed to have a binary membership to an F-formation and to our knowledge, no work has considered refining and enriching this model to label individuals who are partially involved in it. Second, global parameters for the frustum of attention of each person have been used for the entire visual scene. However, social psychology theory has cited the relaxation of the geometric model of an F-formation when consider-

ing the spatial constraints of a room and the furniture in it [8]. Finally, aside from Hung et al. [6], we believe that no other works have seriously addressed the inherently subjective nature of F-formation detection. Our experiments show that by considering the subjectivity of the task, we are better able to model the social scene. That is, by performing associate detection, we show that we can also significantly improve performance on the F-formation detection task.

Concretely, we offer the following contributions; First, we address the novel task of detecting associates of F-formations and propose a novel feature representation that copes with learning from sparse training data. We also show that the state-of-the-art model for full members of F-formations [19] are not appropriate for the modelling of associate behaviour. Second, we model the spatial context of a scene for better F-formation and associate detection by learning a location-dependent frustum of attention of individuals in the scene. Moreover, we address the problem of learning the relative weighting between proximity and orientation given the spatial context of furniture. Third, we contribute new multi-annotator labels on the publicly available Idiap Poster Dataset [6] for modeling associates. Finally, we carry out a deep evaluation and analysis of associates to investigate the complexity of this novel task.

2. Definitions

F-formations and their Associates The psychologist Kendon [8] defined a single conversing group as an F-formation; as a spatial and orientational organization of individuals where each member has equal access to all other members of the group. An F-formation usually consists of three parts, see Figure 1 (a). The o-space is a convex empty space surrounded by the F-formation members, in which every participant orientates themselves inwards, and no external people are allowed. The participants themselves stand in the p-space, which is a narrow strip surrounding the o-space, while the area beyond is called the r-space. Its definition has made it a popular detection task as it relates well to finding maximal cliques in edge-weighted graphs [6, 19, 20]. In practice, a geometric model of a conversing group should be adapted when considering the spatial constraints of a room and the furniture in it [8]. For instance, people talking in front of a laptop may stand closer and look at the same direction (see Figure 1(c)), which maintains an F-formation although their o-space could be violated.

Unlike full members of F-formations, Kendon [8] defines associates to be people who are attached to an F-formation but who are not fully involved in the conversation. Associates can be people who try to join an F-formation but are not fully accepted by the group, or can leave an F-formation abruptly without disturbing the conversation. We name these out-group and in-group associates respectively as the former tends to stand in the r-space while the latter tends to stand in the p-space. Another example

of an associate could be someone who is waiting for a full member (e.g. their spouse) to leave the F-formation and is not interested in engaging in the conversation.

While F-formations can easily be modelled by either maximal cliques [6, 19, 20] or a joint centre-of-focus in the o-space [2], associate behaviours are not so clearly linked to a single set of social cues. Therefore, the associate detection problem requires us to bridge a wider gap and the nature of the problem and how to solve it cannot be so easily translated into a single set of geometric constraints. From the perspective of semantic labelling of a scene, we must also consider that distinguishing full members of F-formations from associates and also singletons is quite important conceptually. Singletons have no social influence on the groups around them. Full F-formation members have the most potential to influence other members of the groups. Meanwhile, associates have the least potential to influence full members but could be influenced by them. Crucially, in-group associates could be mistaken for full F-formation members and out-group associates for singletons.

Frustum of Attention The frustum of attention [19] (or transactional segment, as defined by Kendon [8]) can be considered as a cone-like region extending from the body that represents the spatial and angular extent at which someone is able to see, hear, and potentially touch something or someone else. It represents a three-dimensional space around the human body in which most of our senses and actions are able to be deployed for social interaction. Prior studies have shown that head pose [15, 16, 19], body pose [6], gaze [16, 7], and proximity [6] often provide reliable features for F-formation modeling.

Recent state-of-the-art approaches have tended to use sampling methods to approximate the frustum of attention where the parameters are set carefully by grid search on the entire dataset and the same global model for the frustum of attention is used [19, 2, 13]. There are two main drawbacks of this approach. First, the parameters are likely to overfit on a certain dataset due to the same data being used for training and testing. Second, the variation in F-formation shape caused by the furniture arrangement and non-uniform densities in the crowding of the scene cannot be captured. For example, people can tend to crowd more densely around the area of a bar area even if they are not trying to order drinks or lean on it.

3. Related Work

Exploiting the frustum of attention is very important for detecting F-formations, studies have showed that head pose [15, 16, 19], body pose [5], gaze [16, 7], and proximity [6] often provide reliable patterns. In [22], F-formations are detected by estimating people's position and lower body orientation using only their head position and orientation from a single camera. The modularity cut algorithm [9]

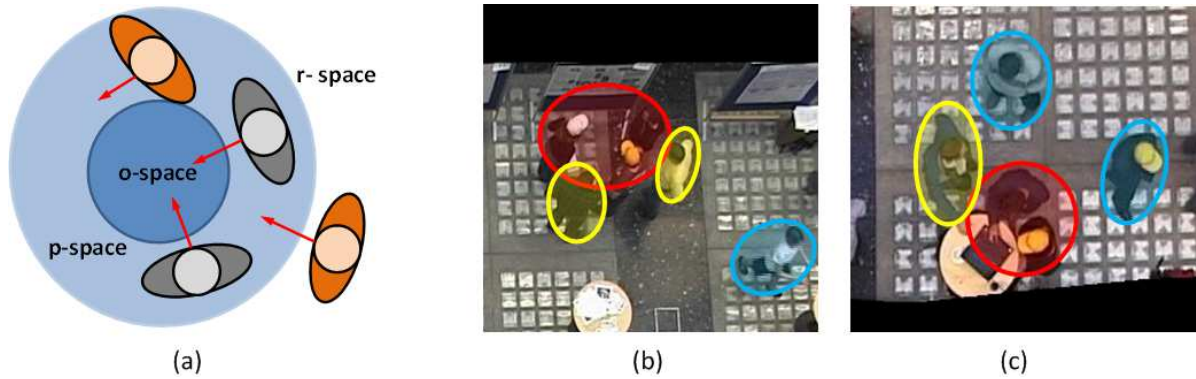


Figure 1. Illustrations of F-formations. (a) The F-formation spaces, gray people stand in the p-space. Red arrows indicate body orientation. Orange people are associates of the F-formation. (b) and (c) example snapshots: F-formations members, associates, and singletons are circled in red, yellow, and blue respectively according to one of our annotators.

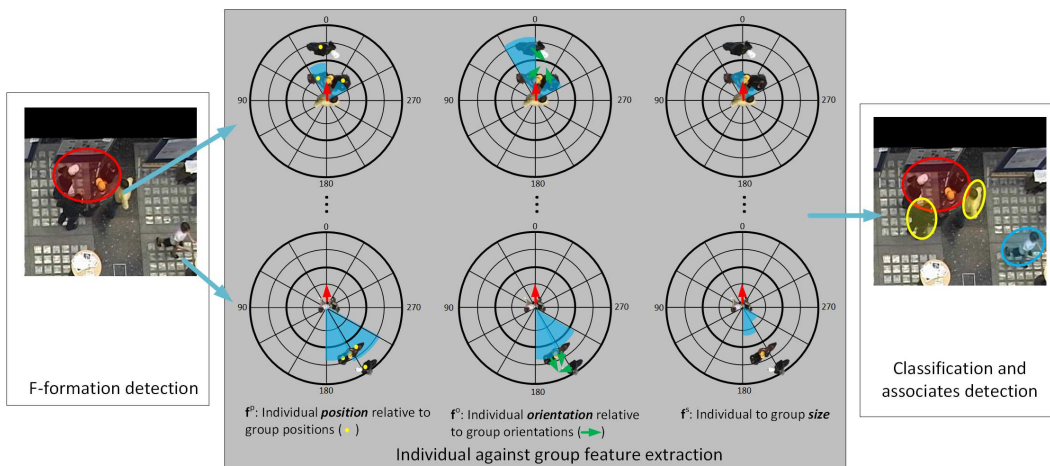


Figure 2. Flow diagram showing the stages of F-formation and associate detection.

was proposed to identify F-formations from automatically extracted trajectories by [23]. To our knowledge, in terms of the treatment of hierarchy in groups, the work of [23] is quite close to ours as they proposed to use eigendecomposition to find centrality in a large mingling group of people. Unfortunately, the data they used was staged but showed participants with high centrality to be those who mingled with more different people.

A Hough voting strategy was proposed in [2], which estimates the locations of o-spaces by density estimation. The size of F-formation was taken into account using a multi-scale Hough voting strategy in [13]. In [6, 19], detecting F-formations is considered as a clustering problem, where each person is defined as a node in the graph, and each edge is the "closeness" between a pair of people. The goal is to find a dominant set [10] in the graph and the edges of the graph are computed based on body orientation and proximity. In [19], the temporal information is added in the dominant set based approach. A density-based approach was proposed in [4] where the final purpose of the task was

to dynamically select camera angles for automated event recording. In [17], temporal patterns of activities were subsequently analyzed. In this paper, we follow the dominant set framework because it gives reliably good results in general [19] and enables a systematic explanation of the learned model so we can interpret better the social phenomena at play in the experimental data. In contrast to the growing numbers of works on F-formation detection, to our knowledge, no one has attempted to detect associates before.

4. Data

We used the publicly available Idiap Poster Data [6]¹, which consists of 3 hours of aerial video of over 50 people during a scientific poster session and coffee break. In this poster session, posters are put around the perimeter of the scene, two small round tables are located in the middle and bottom of the image, a drinks table is located in the bottom right of the image, two entrances are located at the far left and top right of the scene. A screen shot is shown in the left

¹<https://www.idiap.ch/dataset/idiap-poster-data>

of Figure 5. In total, 82 images including 1700 instances of people were annotated by 24 paid annotators, where each image was annotated by 3 annotators. No consecutively selected images contained the same set of formations. We used the positions and body orientation provided separately by Hung et al. [6]. We augmented this data by adding annotations of associates of the F-formations.

We analyzed the annotations to see whether there was full agreement between the annotators about all members of an F-formation and associates. 211 instances of associates were annotated. 84 associates were identified with majority agreement (39.8%) and 34 for full agreement (16%). We computed the F1 score considering one annotation as ground truth and one other annotation as detection for each set of data annotated by the same 3 annotators. The mean and standard deviation of the F1 score are 44% and 13% respectively, which shows that associates are not as straight forward to label compared to F-formations (94.74% mean average F-measure when computing the agreement for F-formations from the data). We consider all the annotated associates have different levels of involvement to groups, that can be perceived by annotators.

To explore the relative angle and orientation relationship between different types of associates of F-formations, we computed histograms of both the distance to, and the relative orientation differences between, an associate and his closest F-formation member as shown in the top and bottom of Figure 5(b) on p. 8 respectively. The relative orientation of associates to their closest F-formation member has a peak in probability mass at 0, and $\pi/3$ while there is only a single peak in the lower histogram. This shows that associates tend to stand similarly closely to their nearest F-formation member. The double peak seen in the relative orientation highlights the idea of two types of associates; those who stand in the p-space of an F-formation but appear less involved in the conversation (in-group associates) and those that stand in the r-space, facing towards the F-formation (out-group associates).

5. Methodology

We detect an associate by modeling its social prior with its associated conversational group (F-formation) based on non-verbal cues where a set of scale (group size) and orientation invariant features are used to train the social prior. The flowchart of the methodology is shown in Figure 2. Given the position and body orientation on the group plane of a set of people, a group detector is first applied to find the conversational groups location (F-formation will be used in the following sections to indicate conversational groups); social prior features are extracted next from every individual; trained classifiers will be used to determine the involvement of a certain people to a F-formation, for instance, F-formation members, associates, or singletons. The modules are described in the following subsections separately.

5.1. Modeling the F-formation as a Dominant set

Building on prior work [6, 19], we exploit the dominant set framework. In an image, people can be represented as a graph $G = (V, E, A)$, where the nodes V are people, E is the set of connections between people, and $A = \{a_{ij}\}$, $i, j \in V$ is an affinity function which defines the "closeness" between each pair of people. Given a subset S of the set of nodes in the graph, the *average weighted degree* of a node $i \in S$ with respect to set S is $k_S(i) = \frac{1}{|S|} \sum_{j \in S, j \neq i} a_{ij}$. The *relative affinity* between node $j \notin S$ and i is $\phi_S(i, j) = a_{ij} - k_S(i)$, and the weight of each i with respect to a set $S = R \cup \{i\}$ is defined as

$$w_S(i) = \begin{cases} 1 & |S|=1 \\ \sum_{j \in R} \phi_R(j, i) w_R(j) & \text{otherwise} \end{cases}, \quad (1)$$

which measures the overall relative affinity between i and the rest of the nodes in S . The relationship between internal and external nodes of a dominant set S is defined as

$$w_S(i) > 0, \quad \forall i \in S \quad (2)$$

$$w_{S \cup \{i\}}(i) < 0, \quad \forall i \notin S. \quad (3)$$

Detecting a dominant set is identical to solving the following standard quadratic programme $\max_{\mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x}$, *s.t.* $\mathbf{x} \in \Delta$, where $\Delta = \{\mathbf{x} \in \mathbb{R}^{|V|} : \sum_{i \in V} x_i = 1, x_i \geq 0, i = 1, \dots, |V|\}$. The indexes of non-zero x_i should correspond to the an F-formation, in such a way that a F-formation can be identified. This optimization problem can be solved with a method from evolutionary game theory, called replicator dynamics. The first-order replicator can be represented as $x_i = x_i \frac{(\mathbf{A}\mathbf{x})_i}{\mathbf{x}^T \mathbf{A} \mathbf{x}}$. Once \mathbf{x} converges, one set of F-formation members are detected. A peeling method is used where the detected group is removed and the replicator dynamics is repeated to find the next F-formation. This peeling method is repeated until the minimum distance of pairwise F-formation members is larger than the maximum distance of detected pairwise F-formation members for a given image. Similar to [6] this enables a stopping criterion that is sensitive to the global context of the scene. For more details, see [6, 10].

5.2. Social involvement features

As described in Section 1 associates have a complex behaviour that is strongly related to the F-formation that they are associated with. They can exist in either the p-space or r-space. Moreover, unlike the maximal clique constraint of full members of F-formations, associates should be mathematically defined with respect to the spatial arrangement of a candidate set of full members of an F-formation. Searching the space of all possible solutions for an associate and F-formation is NP. Fortunately, in practice, associates tend

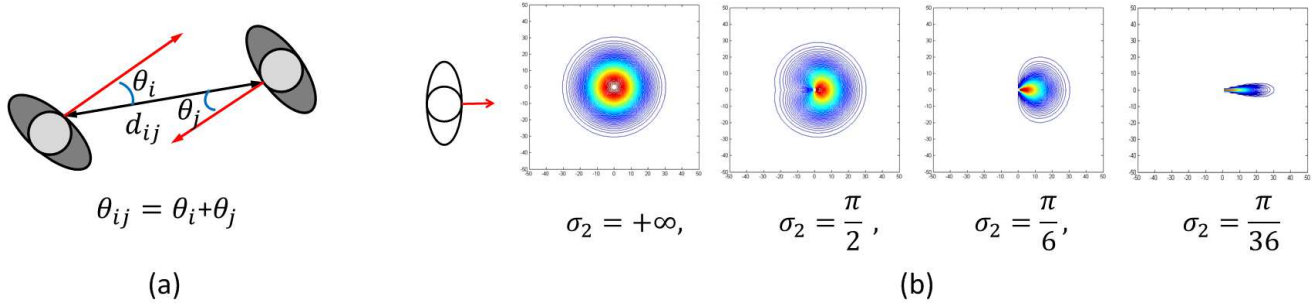


Figure 3. Frustum of attention modeling with body orientation and proximity. (a) Calculation of relative orientation and proximity, (b) frustum of attention map with different parameters. The smaller the σ_2 is, the narrower of frustum attention of a person is.

to be scattered sparsely enough amongst the F-formations in a scene so that the maximal clique assumption for a single F-formation is not severely disrupted by their presence. Therefore in the first instance, using any existing F-formation detection method to reduce the space of possible hypothesis associate and F-formation pairs is reasonable.

Despite this simplification, another challenge still remains. Due to its sparsity, it is unlikely that a sufficient set of examples exist to account for all possible spatial configurations of an associate and F-formation. Therefore, applying similar features that were used to define full members will lead to a representation that is too sparse to learn from. To make sufficiently descriptive features, we hypothesise therefore that they must be both invariant to the rotation of the associate relative to the group, and also insensitive to the size of the group.

To better understand associates and avoid incorrect F-formation detection in the earlier step (*e.g.*, detecting associates as full F-formation members), every individual in the data is considered as an associate candidate, so an associate candidate could be an F-formation member, an associate, or a singleton in reality. Three sets of social prior features $\mathbf{f} = [\mathbf{f}^p, \mathbf{f}^o, \mathbf{f}^s]$, centered at the associate candidate, are extracted to represent the geometric relationship of an associate candidate and its associated F-formation, where the features are based on proximity, body orientation, and group size, respectively. The closest F-formation C to a certain associate candidate \mathbf{p}_a is considered as the associated F-formation of this associate candidate, and \mathbf{p}_k indicates the location of the k^{th} F-formation member in C .

Each set of social prior feature \mathbf{f} is a 12-bin histogram, which is defined based on the angle of the vector between F-formation member \mathbf{p}_k and an associate candidate $\angle(\mathbf{p}_k - \mathbf{p}_a)$, so that every bin covers an angle of $\pi/6$. We define the

m^{th} bin of the three sets of features as

$$\mathbf{f}_m^p = \frac{1}{Z_d \cdot |C_m|} \sum_{k \in C_m} \|\mathbf{p}_k - \mathbf{p}_a\|, \quad (4)$$

$$\mathbf{f}_m^o = \frac{1}{Z_o \cdot |C_m|} \sum_{k \in C_m} (\angle \mathbf{p}_k - \angle \mathbf{p}_a), \quad (5)$$

$$\mathbf{f}_m^s = \frac{1}{Z_s} |C_m|, \quad (6)$$

where the set of F-formation members located in this bin is C_m . We use \mathbf{f}_m^p to represent the average distance between F-formation members in C_m and \mathbf{p}_a , \mathbf{f}_m^o to represent the average relative body orientation between F-formation members in C_m and \mathbf{p}_a , and \mathbf{f}_m^s to represent the relative person density in C_m . The features are normalized by Z_d , Z_o , and Z_s , where Z_d is the maximum proximity between associated F-formation members and associate candidate, $Z_o = 2\pi$, and Z_s is the maximum F-formation size. The middle image in Figure 2 shows examples of the scale or orientation invariant feature representations of an associate and a singleton, which encode people’s relative location, orientation and group size.

Associates detection is challenging because they are likely to be detected as full F-formation members compare to singletons who are usually far away from an F-formation. We use a one-vs-the rest strategy to train an associates detector. In the experiment, we compare a set of classifiers, such as Parzen, RBF SVM, Random Forests, and AdaBoost, with 10 fold cross validation. Parzen classifier gave the best performance on our dataset. In our experiment, we used 211 instances of annotated associates, 235 full-agreement singletons and 450 full-agreement F-formations as training data.

5.3. Training the affinity matrix

To detect F-formations in a complex environment, we need to model the variation of the density of geometric variations of potential F-formations in the space. To capture this variation, the affinity matrix \mathbf{A} is key. In this paper,

we only consider the proximity and body orientation. The "closeness" between people i and j is defined as

$$a_{ij} = e^{-\frac{d_{ij}^2}{\sigma_1^2} - \frac{\theta_{ij}^2}{\sigma_2^2}}, \quad (7)$$

where d_{ij} is the Euclidean distance between two people, θ_{ij} is the sum of difference between each body orientation and the angle of the vector between two people (see Figure 3), and σ_1 and σ_2 are the parameters to be learned. As the values of σ_1 and σ_2 decrease, a person is likely to stand closer and angle more directly towards the others in the F-formation (see Figure 5 (a)). Likewise, as σ_1 and σ_2 increase, members of an F-formation will tend to stand further apart and orientate themselves less directly towards others (see Figure 5 (a)). The objective function is defined as

$$\ell = \sum_{n=1}^N 1 - \frac{C^{\{n\}} \cap \hat{C}^{\{n\}}}{C^{\{n\}} \cup \hat{C}^{\{n\}}} \quad (8)$$

where n is the index of an F-formation in an image, N is the total number of annotated F-formations, and $C^{\{n\}}$ and $\hat{C}^{\{n\}}$ are the n th detected set of F-formation members and its corresponding annotation respectively. During training, we consider a detection C and an annotation \hat{C} to match with each other if $\frac{|C \cap \hat{C}|}{|C \cup \hat{C}|} \geq \frac{2}{3}$. Considering that the shape of the F-formation can be influenced by the furniture arrangement, we learn parameters σ_1 and σ_2 as a function of a person's location \mathbf{p} . We only update the parameters once per person when the detection goes wrong in a passive-aggressive way [1].

$$\sigma_s(\mathbf{p}) = \sigma_s(\mathbf{p}) - g_s(C)\Delta\sigma_s, \quad s \in \{1, 2\}. \quad (9)$$

Here, $\Delta\sigma_s$ is the basic step size, which is set to a small value. An adaptive parameter g helps to adapt to different F-formation geometric variations. Given F-formation C , the adaptive parameter g is defined as

$$g_1(C) = y \frac{\|\sum_{i,j \in \hat{C}^{\{n\}}} \hat{d}_{ij} - \sum_{i,j \in C^{\{n\}}} d_{ij}\|}{\sum_{i,j \in \hat{C}^{\{n\}}} \hat{d}_{ij}}, \quad (10)$$

$$g_2(C) = y \frac{\|\sum_{i,j \in \hat{C}^{\{n\}}} \hat{\theta}_{ij} - \sum_{i,j \in C^{\{n\}}} \theta_{ij}\|}{\sum_{i,j \in \hat{C}^{\{n\}}} \hat{\theta}_{ij}}, \quad (11)$$

where $y \in \{-1, 1\}$, $y = 1$ indicates a false negative F-formation member in C , while $y = -1$ indicates a false positive member. Here \hat{d} and $\hat{\theta}$ are the manually annotated proximity and frustum of attention. In each iteration, we update each person's location in the F-formation.

6. Experiment

6.1. Experiment setup

In the experiment, we initialized $\sigma_1 = 40, \sigma_2 = 30$ for training, whose basic update step sizes were set to

$\Delta\sigma_1 = 0.1$ and $\Delta\sigma_2 = \pi/720$ respectively. The number of iterations of training for detecting F-formation and associates were both set to 300. Considering that the training samples in each precise location were not distributed densely over the images, we divided the images into blocks of 45×45 pixels where all people located in the same block shared the same learned parameters. We trained using each of the 3 annotations separately, applying 10 fold cross validation for each. Finally, the position and body orientations used to train our models came from the annotations of the Idiap poster data provided by Hung et al. [6].

For evaluation, we consider a group as correctly estimated if at least $(T \cdot |C|)$ of their members are detected, where $|C|$ is the cardinality of the labeled group C , and $T \in [0, 1]$ is an arbitrary threshold; in [2], the scoring threshold $T = 2/3$, corresponds to finding at least two thirds of the members of a group. Here we also consider $T = 1$, to mean that a group is correctly detected only if all members are labeled correctly. From these metrics we calculate the precision, recall and F1 measures in each frame, averaging them over all the frames and the three sets of annotations. Associates are evaluated by calculating precision, recall and F1 score in the same way, where only the harder $T = 1$ criterion for success is used. Here, a baseline detector global-F is added, which only uses the initialized training value $\sigma_1 = 40, \sigma_2 = 30$ for detecting F-formation. We also compared the performance of our spatially-aware F-formation detector (Spatial-F) with state-of-the-art DSFF [6], HFF [2], ACCVKL [19], and ACCVJS [19].

Since we are the first to approach the task of detecting associates, we create three baseline detectors to compare with our proposed associate detector (social-A). Each baseline result was generated using the annotated data and not detections. First, **SA** labels all people who are not in an F-formation (mostly singletons) as associates. Second, **RA** labels people as associates of an F-formation if their distance to it is less than or equal to the average distance between pairwise members of F-formations according to the entire labeled data. Third, **ADA** is set based on the average disagreement between annotators where for each pair, we treated one annotation as a detected result to compute performance against another annotation. We also compared performances with different feature combinations (\mathbf{p} : proximity features, \mathbf{o} : orientation features, and \mathbf{s} : group size features). The associates detector global-A extracts features based on global-F F-formation detection.

Finally, we analysed how associate detection can help improve F-formation detection. As the F-formation detector has problems mostly with in-group associates, we used the detected associates to clean up false positives in a detected F-formation. The performance of Spatial-F and global-F was evaluated with the $T = 1$ hard criterion using F-formations annotated with full agreement.

Table 1. F-formation detection results with soft ($T = 2/3$) and hard ($T = 1$) criteria for deciding on whether an F-formation is correctly detected.

Method	T=2/3			T=1		
	Prec.	Rec.	F1	Prec.	Rec.	F1
DSFF [6]	0.93	0.92	0.92	0.81	0.81	0.81
HFF [2]	0.93	0.96	0.94	0.81	0.84	0.83
ACCVKL [19]	0.90	0.94	0.92	-	-	-
ACCVJS [19]	0.92	0.96	0.94	-	-	-
global-F	0.87	0.92	0.89	0.72	0.76	0.74
spatial-F	0.91	0.98	0.94	0.91	0.98	0.94

6.2. F-formation Detection Results

Two examples of the learned values for σ_1 and σ_2 with respect to the spatial context, are shown in Figure 5 (a). People in the top F-formation standing side-by-side tend to have a large σ_2 , while people in the bottom F-formation standing face-to-face tend to have a small σ_2 .

From Table 1, for $T = 2/3$, our detector (spatial-F) shows competitive performance to the state-of-art. This is because tuning a global value of σ can already produce a good approximation of the clean F-formation shape, particularly as the soft detection threshold already considers partially detected members of an F-formation to be sufficient, enabling a softening of the need for strongly circular formations. However, when considering the harsher criterion $T = 1$, our detector (spatial-F) significantly out-performs the state-of-the-art, even with a cross-validated comparison. We can also see that the spatial-F detector performs equally good with both criteria ($T = 2/3$ or 1), which shows the accuracy of our detector is very high.

6.3. Results of Detecting Associates of F-formations

Table 2 shows that our proposed associate detector (social-A) significantly outperforms the three baselines (SA, RA and ADA), which means there are indeed certain patterns of associate behaviour that differs from the behaviour of singletons. We can also see from the performance ADA that it is also difficult for people to agree on who associates are. It also shows that social-A(p+o) with only proximity and orientation features can almost achieve the performance of the complete set of features. Interestingly, global-A shows features extracted with a less accurate F-formation detector can still obtain a similar performance with social-A where a more accurate F-formation detector spatial-F was used. This can be explained as our feature represents prototype-like F-formation structures, which can tolerate certain errors on less perfect F-formation detections.

To understand more about associates, some examples of them are shown in Figure 4. The red dots indicate the members' positions in an F-formation, the small red lines indicate everyone's orientation, the yellow dots indicate the

Table 2. Associate detection results. SA: labels all singletons as associates, RA: labels people close to F-formation as associates, UA: performance based on annotator disagreement, global-A: use global-F detector to extract features, and social-A: our proposed detector (details in Sec. 6.1).

Method	Prec.	Rec.	F1
SA	0.06	1.00	0.11
RA	0.11	0.84	0.19
ADA	0.44	0.44	0.44
global-A(p+o+s)	0.89	0.59	0.71
social-A(p)	0.87	0.58	0.69
social-A(o)	0.91	0.55	0.69
social-A(s)	0.78	0.53	0.63
social-A(p+o)	0.89	0.57	0.70
social-A(p+s)	0.85	0.56	0.67
social-A(o+s)	0.91	0.56	0.69
social-A(p+o+s)	0.89	0.59	0.71

Table 3. F-formation detection with associate detection feedback, results are evaluated only on F-formations annotated with full-agreement. FB-global-F and FB-spatial-F are detectors with associate detection feedback (details in Sec. 6.1).

Method	Prec.	Rec.	F1
global-F	0.75	0.94	0.83
FB-global-F	0.82	0.94	0.88
spatial-F	0.76	1.00	0.86
FB-spatial-F	0.84	1.00	0.91

correctly detected associates, the blue dots are correctly detected singletons, and the green dots show associates that were missed by the detector. From left to right, the first two images show that our detector can successfully detect associates who are in the r-space (See Figure 1(a)) trying to join an F-formation but who are not accepted by its members. The third and fourth images show that our detector can detect associates who are still in the F-formation p-space but not fully involved in the group. This conforms our analysis of the orientation and proximity of associates in Section 4 Figure 5(b).

We simulated tracking drifts on the manual labels of position and body orientation to compare the robustness of our method spatial-F with global-F on noisy test data. Figure 4 (b) shows that our detector spatial-F in general performs better than the detector with global parameters global-F, however, our detector can tolerate less noise by looking at the decay rate because our learned parameters are sensitive to the location changing. As a person width is approximately 20 pixels in the image, the performance of our method starts to drop faster when the deviation of Gaussian noise is around half person width. It means our method should perform well using a reasonably robust visual tracker.

From Table 3, we can see that using the feedback of the detected associates, false positive F-formation members are



Figure 4. (a): example associate detection results: Red dots - members of an F-formation; red lines - body orientation; yellow dots - correctly detected associates; blue dots - correctly detected singletons; and green dots - missed associate detections. (b): F1 score of F-formation detectors spatial-F and global-F and associates detectors social-A and global-A with noisy test data.

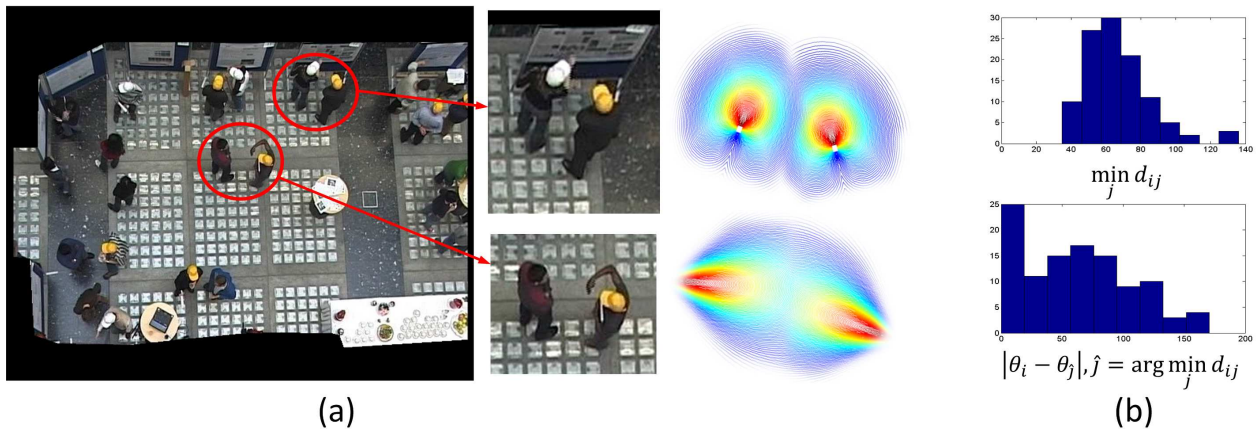


Figure 5. (a): learned frustum of attention in two cases. (b): histograms of both relative orientation differences between an associate and also distance to closest nearest F-formation member.

removed, so that the precisions are improved significantly.

7. Conclusion

In this paper, we addressed the novel task of detecting associates of F-formations. We introduced a novel multi-annotator annotations for associates of F-formations, and two methods for detecting them. Using our model, we were also able to discover patterns in proximity and orientation in the behaviours of associates that enable significant improvement over baseline methods with a detection rate of 71% F-measure. We proposed a spatial-context-aware F-formation detector, which models people’s frustum of attention in a principled way while considering the influence of the social and spatial context. The method is in general more adaptive to different datasets so for example, different frustum of attention parameters can be learned from scenarios with a non-uniform density of crowding. Our proposed method showed competitive performance, even when training the model parameters on less data.

By cleaning the detected in-group associates from the detected F-formations, we were also able to significantly improve F-formation detection performance in all cases where there was full agreement amongst annotators on the full-members of each F-formation. Surprisingly, although learning a spatial-context specific frustum of attention led to better F-formation detection, when using the output of this models to detect associates, the performance for associate detection was not better than when F-formations were detected with a spatial-context free frustum parameters.

In summary, to our knowledge, this constitutes the first attempt on the challenging problem of automatically estimating conversational involvement levels in visual scenes of mingling.

Acknowledgments. This work has partly been supported by the European Commission under contract number FP7-ICT-600877 (SPENCER). The authors thank Jan van Gemert and Julian Kooij for helpful discussions.

References

- [1] K. Crammer, O. Dekel, J. Keshet, S. Shalev-Shwartz, and Y. Singer. Online passive-aggressive algorithms. *Journal of Machine Learning Research*, 7(Mar):551–585, 2006.
- [2] M. Cristani, L. Bazzani, G. Paggetti, A. Fossati, D. Tosato, A. Del Bue, G. Menegaz, and V. Murino. Social interaction discovery by statistical analysis of F-formations. In *BMVC*, pages 1–12, 2011.
- [3] M. Cristani, R. Raghavendra, A. D. Bue, and V. Murino. Human behavior analysis in video surveillance: A social signal processing perspective. *Neurocomputing*, 100(0):86–97, 2013. Special issue: Behaviours in video.
- [4] T. Gan, Y. Wong, D. Zhang, and M. S. Kankanhalli. Temporal encoded f-formation system for social interaction detection. In *Proceedings of the 21st ACM international conference on Multimedia*, pages 937–946. ACM, 2013.
- [5] G. Groh, A. Lehmann, J. Reimers, M. R. Frieß, and L. Schwarz. Detecting social situations from interaction geometry. In *Social Computing (Social-Com), 2010 IEEE Second International Conference on*, pages 1–8. IEEE, 2010.
- [6] H. Hung and B. Kröse. Detecting f-formations as dominant sets. In *Proceedings of the 13th international conference on multimodal interfaces*, pages 231–238. ACM, 2011.
- [7] N. Jovanović et al. Towards automatic addressee identification in multi-party dialogues. Association for Computational Linguistics, 2004.
- [8] A. Kendon. *Conducting interaction: Patterns of behavior in focused encounters*, volume 7. CUP Archive, 1990.
- [9] M. E. Newman. Modularity and community structure in networks. *Proceedings of the National Academy of Sciences*, 103(23):8577–8582, 2006.
- [10] M. Pavan and M. Pelillo. Dominant sets and pairwise clustering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(1):167–172, 2007.
- [11] L. Porzi, S. Rota Bulò, B. Lepri, and E. Ricci. Predicting and understanding urban perception with convolutional neural networks. In *Proceedings of the 23rd Annual ACM Conference on Multimedia Conference, MM ’15*, pages 139–148, New York, NY, USA, 2015. ACM.
- [12] D. Santani and D. Gatica-Perez. Loud and trendy: Crowdsourcing impressions of social ambiance in popular indoor urban places. In *Proceedings of the 23rd Annual ACM Conference on Multimedia Conference, MM ’15*, pages 211–220, New York, NY, USA, 2015. ACM.
- [13] F. Setti, O. Lanz, R. Ferrario, V. Murino, and M. Cristani. Multi-scale f-formation discovery for group detection. In *ICIP*, pages 3547–3551, 2013.
- [14] A. W. M. Smeulders, M. Worring, S. Santini, A. Gupta, and R. Jain. Content-based image retrieval at the end of the early years. *IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE*, 22(12):1349–1380, 2000.
- [15] K. Smith, S. O. Ba, J.-M. Odobez, and D. Gatica-Perez. Tracking the visual focus of attention for a varying number of wandering people. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 30(7):1212–1229, 2008.
- [16] R. Subramanian, J. Staiano, K. Kalimeri, N. Sebe, and F. Pianesi. Putting the pieces together: multimodal analysis of social attention in meetings. In *Proceedings of the international conference on Multimedia*, pages 659–662. ACM, 2010.
- [17] K. Tran, A. Gala, I. Kakadiaris, and S. Shah. Activity analysis in crowded environments using social cues for group discovery and human interaction modeling. *Pattern Recognition Letters*, 44:49–57, 2014.
- [18] R. Triebel, K. Arras, R. Alami, L. Beyer, S. Breuers, R. Chatila, M. Chetouani, D. Cremers, V. Evers, M. Fiore, H. Hung, O. Islas Ramirez, M. Joosse, H. Khambhaita, T. Kucner, B. Leibe, A. Lilienthal, T. Linder, M. Lohse, M. Magnusson, B. Okal, L. Palmieri, U. Rafi, M. van Rooij, and L. Zhang. Spencer: A socially aware service robot for passenger guidance and help in busy airports. In *Conference on Field and Service Robotics (FSR)*, 2015.
- [19] S. Vascon, E. Z. Mequanint, M. Cristani, H. Hung, M. Pelillo, and V. Murino. A Game-Theoretic Probabilistic Approach for Detecting Conversational Groups. In *ACCV*, 2014.
- [20] S. Vascon, E. Z. Mequanint, M. Cristani, H. Hung, M. Pelillo, and V. Murino. Detecting conversational groups in images and sequences: A robust game-theoretic approach. *Computer Vision and Image Understanding*, pages –, 2015.
- [21] A. Vinciarelli, M. Pantic, D. Heylen, C. Pelachaud, I. Poggi, F. D’Errico, and M. Schroeder. Bridging the gap between social animal and unsocial machine: A survey of social signal processing. *IEEE Transactions on Affective Computing*, 3(1):69–87, 2012.
- [22] N. Yasuda, K. Kakusho, T. Okadome, T. Funatomi, and M. Iiyama. Recognizing conversation groups in

an open space by estimating placement of lower bodies. In *Systems, Man and Cybernetics (SMC), 2014 IEEE International Conference on*, pages 544–550, Oct 2014.

- [23] T. Yu, S. Lim, K. A. Patwardhan, and N. Krahnstoeber. Monitoring, Recognizing and Discovering Social Networks. In *CVPR*, 2009.