

PERI: Part Aware Emotion Recognition In The Wild

Akshita Mittel¹ and Shashank Tripathi²

¹ NVIDIA

amittel@nvidia.com

² Max Planck Institute for Intelligent Systems, Tübingen, Germany

stripathi@tue.mpg.de

Abstract. Emotion recognition aims to interpret the emotional states of a person based on various inputs including audio, visual, and textual cues. This paper focuses on emotion recognition using visual features. To leverage the correlation between facial expression and the emotional state of a person, pioneering methods rely primarily on facial features. However, facial features are often unreliable in natural unconstrained scenarios, such as in crowded scenes, as the face lacks pixel resolution and contains artifacts due to occlusion and blur. To address this, methods focusing on *in the wild* emotion recognition exploit full-body person crops as well as the surrounding scene context. While effective, in a bid to use body pose for emotion recognition, such methods fail to realize the potential that facial expressions, when available, offer. Thus, the aim of this paper is two-fold. First, we demonstrate a method, PERI, to leverage both body pose and facial landmarks. We create *part aware spatial* (PAS) images by extracting key regions from the input image using a mask generated from both body pose and facial landmarks. This allows us to exploit body pose in addition to facial context whenever available. Second, to reason from the PAS images, we introduce context infusion (Cont-In) blocks. These blocks attend to part-specific information, and pass them onto the intermediate features of an emotion recognition network. Our approach is conceptually simple and can be applied to any existing emotion recognition method. We provide our results on the publicly available in the wild EMOTIC dataset. Compared to existing methods, PERI achieves superior performance and leads to significant improvements in the mAP of emotion categories, while decreasing Valence, Arousal and Dominance errors. Importantly, we observe that our method improves performance in both images with fully visible faces as well as in images with occluded or blurred faces.

1 Introduction

The objective of emotion recognition is to recognise how people feel. Humans function on a daily basis by interpreting social cues from around them. Lecturers can sense confusion in the class, comedians can sense engagement in their audience, and psychiatrists can sense complex emotional states in their patients.

As machines become an integral part of our lives, it is imperative that they understand social cues in order to assist us better. By making machines more aware of context, body language, and facial expressions, we enable them to play a key role in numerous situations. This includes monitoring critical patients in hospitals, helping psychologists monitor patients they are consulting, detecting engagement in students, analysing fatigue in truck drivers, to name a few. Thus, emotion recognition and social AI have the potential to drive key technological advancements in the future.

Facial expressions are one of the biggest indicators of how a person feels. Therefore, early work in recognizing emotions focused on detecting and analyzing faces [2, 12, 33, 37]. Although rapid strides have been made in this direction, such methods assume availability of well aligned, fully visible and high-resolution face crops [8, 18, 21, 29, 31, 35, 39]. Unfortunately, this assumption does not hold in realistic and unconstrained scenarios such as internet images, crowded scenes, and autonomous driving. In the wild emotion recognition, thus, presents a significant challenge for these methods as face crops tend to be low-resolution, blurred or partially visible due to factors such as subject’s distance from the camera, person and camera motion, crowding, person-object occlusion, frame occlusion etc. In this paper, we address in-the-wild emotion recognition by leveraging face, body and scene context in a robust and efficient framework called Part-aware Emotion Recognition In the wild, or PERI.

Research in psychology and affective computing has shown that body pose offers significant cues on how a person feels [5, 7, 20]. For example, when people are interested in something, they tilt their head forward. When someone is confident, they tend to square their shoulders. Recent methods recognize the importance of body pose for emotion recognition and tackle issues such as facial occlusion and blurring by processing image crops of the entire body [1, 19, 22, 24, 38, 41]. Kosti et al. [22, 24] expand upon previous work by adding scene context in the mix, noting that the surrounding scene plays a key role in deciphering the emotional state of an individual. An illustrative example of this could be of a person crying at a celebration such as graduation as opposed to a person at a funeral. Both individuals can have identical posture but may feel vastly different set of emotions. Huang et al. [19] expanded on this by improving emotion recognition using body pose estimations.

In a bid to exploit full body crops, body keypoints and scene context, such methods tend to ignore part-specific information such as shoulder position, head tilt, facial expressions, etc. which, when available, serve as powerful indicators of the emotional state. While previous approaches focus on either body pose or facial expression, we hypothesize that a flexible architecture capable of leveraging both body and facial features is needed. Such an architecture should be robust to lack of reliable features on both occluded/blurred body or face, attend to relevant body parts and be extensible enough to include context from the scene. To this end, we present a novel representation, called part-aware spatial (PAS) images that encodes both facial and part specific features and retains pixel-alignment relative to the input image. Given a person crop, we generate a part-aware mask

by fitting Gaussian functions to the detected face and body landmarks. Each Gaussian in the part-aware mask represents the spatial context around body and face regions and specifies key regions in the image the network should attend to. We apply the part-aware mask on the input image which gives us the final PAS image (see Fig. 2). The PAS images are agnostic to occlusion and blur and take into account both body and face features.

To reason from PAS images, we propose novel context-infusion (Cont-In) blocks to inject part-aware features at multiple depths in a deep feature backbone network. Since the PAS images are pixel-aligned, each Cont-In block implements explicit attention on part-specific features from the input image. We show that as opposed to *early fusion* (e.g. channel-wise concatenation) of PAS image with input image \mathbf{I} , or *late fusion* (concatenating the features extracted from PAS images just before the final classification), Cont-In blocks effectively utilize part-aware features from the image. Cont-In blocks do not alter the architecture of the base network, thereby allowing Imagenet pretraining on all layers. The Cont-In blocks are designed to be easy to implement, efficient to compute and can be easily integrated with any emotion recognition network with minimal effort.

Closest to our work is the approach of Gunes et al. [15] which combines visual channels from face and upper body gestures for emotion recognition. However, unlike PERI, which takes unconstrained in the wild monocular images as input, their approach takes two high-resolution camera streams, one focusing only on the face and other focusing only on the upper body gestures from the waist up. All of the training data in [15] is recorded in an indoor setting with a uniform background, single subject, consistent lighting, front-facing camera and fully visible face and body; a setting considerably simpler than our goal of emotion recognition in real-world scenarios. Further, our architecture and training scheme is fundamentally different and efficiently captures part-aware features from monocular images.

In summary, we make the following contributions:

1. Our approach, PERI, advances in the wild emotion recognition by introducing a novel representation (called PAS images) which efficiently combines body pose and facial landmarks such that they can supplement one another.
2. We propose context infusion (Cont-In) blocks which modulate intermediate features of a base emotion recognition network, helping in reasoning from both body poses and facial landmarks. Notably, Cont-In blocks are compatible with any existing emotion recognition network with minimal effort.
3. Our approach results in significant improvements compared to existing approaches in the publicly-available in the wild EMOTIC dataset [23]. We show that PERI adds robustness under occlusion, blur and low-resolution input crops.

2 Related Work

Emotion recognition is a field of research with the objective of interpreting a person’s emotions using various cues such as audio, visual, and textual inputs.

Preliminary methods focused on recognising six basic discrete emotions defined by the psychologists Ekman and Friesen [9]. These include anger, surprise, disgust, enjoyment, fear, and sadness. As research progressed, datasets, such as the EMOTIC dataset [22, 23, 24], have expanded on these to provide a wider label set. A second class of emotion recognition methods focus not on the discrete classes but on a continuous set of labels described by Mehrabian [30] including Valence (V), Arousal (A), and Dominance (D). We evaluate the performance of our model using both the 26 discrete classes in the EMOTIC dataset [23], as well as valence, arousal, and dominance errors. Our method works on visual cues, more specifically on images and body crops.

Emotion recognition using facial features. Most existing methods in Computer Vision for emotion recognition focus on facial expression analysis [2, 12, 33, 37]. Initial work in this field was based on using the Facial Action Coding System (FACS) [4, 10, 11, 26, 34] to recognise the basic set of emotions. FACS refers to a set of facial muscle movements that correspond to a displayed emotion, for instance raising the inner eyebrow can be considered as a unit of FACS. These methods first extract facial landmarks from a face, which are then used to create facial action units, a combination of which are used to recognise the emotion. Another class of methods use CNNs to recognize the emotions [2, 19, 22, 23, 32, 42]. For instance, Emotionnet [2] uses face detector to obtain face crops which are then passed into a CNN to get the emotion category. Similar to these methods, we use facial landmarks in our work. However, uniquely, the landmarks are used to create the PAS contextual images, which in turn modulate the main network through a series of convolutions layers in the Cont-In blocks.

Emotion recognition using body poses. Unlike facial emotion recognition, the work on emotion recognition using body poses is relatively new. Research in psychology [3, 13, 14] suggests that cues from body pose, including features such as hip, shoulder, elbow, pelvis, neck, and trunk can provide significant insight into the emotional state of a person. Based on this hypothesis, Crenn et al. [6] sought to classify body expressions by obtaining low-level features from 3D skeleton sequences. They separate the features into three categories: geometric features, motion features, and fourier features. Based on these low-level features, they calculate meta features (mean and variance), which are sent to the classifier to obtain the final expression labels. Huang et al. [40] use a body pose extractor built on Actional-Structural GCN blocks as an input stream to their model. The other streams in their model extract information from images and body crops based on the architecture of Kosti et al. [22, 24]. The output of all the streams are concatenated using a fusion layer before the final classification. Gunes et al. [15] also uses body gestures. Similar to PERI, they use facial features by combining visual channels from face and upper body gestures. However, their approach takes two high-resolution camera streams, one focusing only on the face and other focusing only on the upper body gestures, making them unsuitable for unconstrained settings. For our method, we use two forms of body posture information, body crops and body pose detections. Body crops taken from

the original input image are passed into one stream of our architecture. The intermediate features of the body stream are then modulated at regular intervals using Cont-In blocks, which derive information from the PAS image based on body pose and facial landmarks.

Adding visual context from the entire image. The most successful methods for emotion recognition in the wild use context from the entire image as opposed to just the body or facial crop. Kosti et al. [22, 24] were among the first to explore emotion recognition in the wild using the entire image as context. They introduced the EMOTIC dataset [23] on which they demonstrated the efficacy of a two-stream architecture where one of the streams is supplied with the entire image while the other is supplied with body crops. Gupta et al. [16] also utilise context from the entire image using an image feature extraction stream. Additionally, the facial crops from the original image are passed through three modules; a facial feature extraction stream, an attention block and finally a fusion network. The attention block utilizes the features extracted from the full image to additionally modulate the features of the facial feature extraction stream. However, unlike Kosti et al. they focus on recognising just the valence of the entire scene. Zhang et al. [42] also use context from an image. Their approach uses a Region Proposal Network (RPN) to detect nodes which then form an affective graph. This graph is fed into a Graph Convolutional Network (GCN) similar to Mittal et al. [32]. Similar to Kosti et al. the second CNN stream in their network extracts the body features. Lee et al. [25] present CAERNet, which consists of two subnetworks. CAERNet is two-stream network where one stream works with facial crops and the other in which both facial expression and context (background) are extracted. They use an adaptive fusion network in order to fuse the two streams. Mittal et al. [32] take context a step further. Similar to our approach, they use both body crops and facial landmarks. However, akin to Huang et al. [40] they pass body crops and facial landmarks as a separate stream. Their architecture consists of three streams. In addition to the body pose and facial landmark stream, the second stream extracts information from the entire image where the body crop of the person has been masked out. The third stream adds modality in one of two ways. They first encode spatio-temporal relationships using a GCN network similar to [42], these are then passed through the third stream. The other method uses a CNN which parses depth images in the third stream. Similar to these methods, PERI maintains two streams, where one of the stream extracts meaningful context from the entire image while the other focuses on the individual person.

3 Method

The following section describes our framework to effectively recognize emotions from images in the wild. Facial expressions, where available, provide key insight to the emotional state of a person. We find a way to represent body pose and facial landmarks such that we can utilise both set of features subject to their availability in the input image. Concretely, we use MediaPipe’s Holistic model [28] to obtain

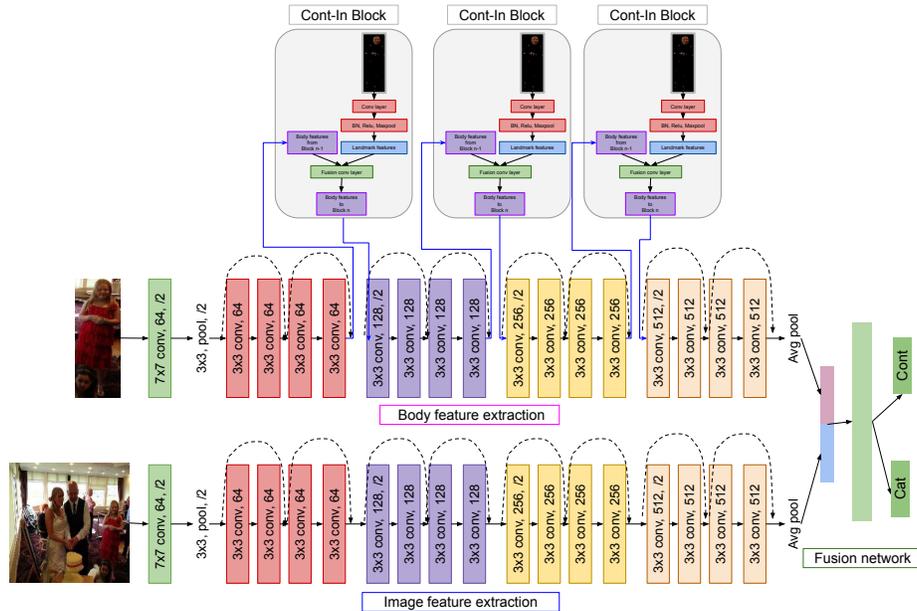


Fig. 1: The overall architecture that consists of a two-stream Resnet-18 network along with the Cont-In Blocks that modulate features after every intermediate Resnet block using the Part Aware Spatial Context images (PAS).

landmarks for face and body. These landmarks are then used to create our part aware spatial (PAS) images. Part-specific context from the PAS images is learnt from our context infusion (Cont-In) blocks which modulate the intermediate features of a emotion detection network. Fig. 1 shows the overall framework that we use for our emotion recognition pipeline. A more detailed view of our Cont-In blocks can be seen in Fig. 2.

3.1 MediaPipe Holistic model

In order to obtain the body poses and facial landmarks, we use the MediaPipe Holistic pipeline [28]. It is a multi-stage pipeline which includes separate models for body pose and facial landmark detection. The body pose estimation model is trained on 224×224 input resolution. However, detecting face and fine-grained facial landmarks requires high resolution inputs. Therefore, the MediaPipe Holistic pipeline first estimates the human pose and then finds the region of interest for the face keypoints detected in the pose output. The region of interest is upsampled and the facial crop is extracted from the original resolution input image and is sent to a separate model for fine-grained facial landmark detection.

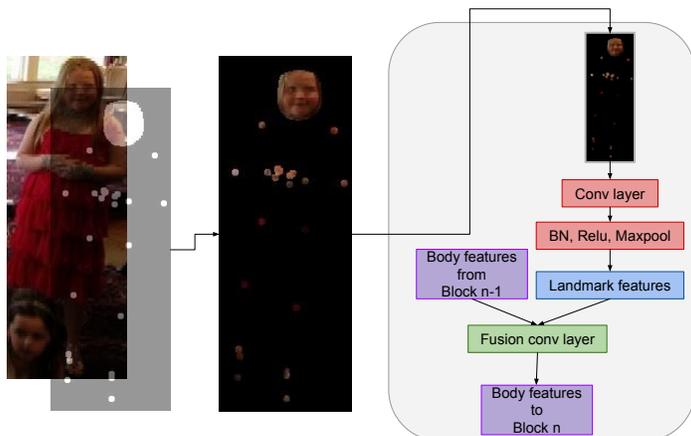


Fig. 2: (Left) An input image along with the mask (\mathbf{B}') created by applying a Gaussian function with $\sigma = 3$. The mask is binarised and used to create the PAS image (\mathbf{P}) in the middle. (Right) Architecture of the Cont-In block that uses the PAS images (\mathbf{P}) to modulate the Resnet features between each intermediate block. Here the input features from the $n - 1^{th}$ Resnet block are passed in and modulated features are passed to the n^{th} block are shown in purple.

3.2 The Emotic Model

The baseline of our paper is the the two-stream CNN architecture from Kosti et. al [22, 24]. The paper defines the task of *emotion recognition in context*, which considers both body pose and scene context for emotion detection. The architecture takes as input the body crop image, which is sent to the body feature extraction stream, and the entire image, which is sent to the image feature extraction stream. The outputs from the two streams are concatenated and combined through linear classification layers. The model outputs classification labels from 26 discrete emotion categories and 3 continuous emotion dimensions, *Valence*, *Arousal* and *Dominance* [30]. The 2 stream architecture is visualized in our pipeline in Fig. 1. In order to demonstrate our idea, we stick with the basic Resnet-18 [17] backbone for both the streams.

3.3 Part aware spatial image

One contribution of our framework is how we combine body pose information along with facial landmarks such that we can leverage both sets of features and allow them to complement each other subject to their availability. In order to do so we have three main stages to our pipeline. First, we use the MediaPipe Holistic model to extract the keypoints as described in Sec. 3.1. Here we get two sets of keypoint coordinates for each body crop image \mathbf{I} . The first set of N coordinates describe the body landmarks \mathbf{b}_i where $i \in (0, N)$. The second set of

M coordinates describe the location of the facial landmarks \mathbf{f}_j where $j \in (0, M)$. For simplicity, we combine all detected landmarks and denote them as \mathbf{b}_k where $k \in (0, M + N)$. We take an all black mask $\mathbf{B} \in \mathbb{R}^{1 \times H \times W}$ the same size as the body crop, and fit a Gaussian kernel to every landmark in the original image as

$$\mathbf{b}'_k = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (1)$$

The part-aware mask $\mathbf{B}' \in \mathbb{R}^{(1 \times H \times W)}$ is created by binarizing \mathbf{b}'_k using a constant threshold ρ , such that

$$\mathbf{B}'(x) = \begin{cases} 1 & \text{if } \|\mathbf{x} - \mathbf{b}_k\|_2 \leq \rho, \\ 0 & \text{if } \|\mathbf{x} - \mathbf{b}_k\|_2 > \rho, \end{cases} \quad (2)$$

where \mathbf{x} is the coordinates of all pixels in \mathbf{B} . The distance threshold ρ is determined empirically.

Finally, to obtain the part aware spatial (PAS) image $\mathbf{P} \in \mathbb{R}^{3 \times H \times W}$, the part-aware mask is applied to the input body crop \mathbf{I} using channel-wise hadamard product,

$$\mathbf{P} = \mathbf{I} \otimes \mathbf{B}' \quad (3)$$

This process can be visualized in Fig. 2 (left).

3.4 Context Infusion Blocks

To extract information from PAS images, we explore *early fusion*, which simply concatenates PAS with the body crop image \mathbf{I} in the body feature extraction stream of our network. We also explore *late fusion*, concatenating feature maps derived from PAS images before the fusion network. However, both of these approaches failed to improve performance. Motivated by the above, we present our second contribution, the Context Infusion Block (Cont-In) which is an architectural block that utilizes the PAS contextual image to condition the base network. We design Cont-In blocks such that they can be easily introduced in any existing emotion recognition network. Fig. 2 shows the architecture of a Cont-In block in detail. In PERI, the body feature extraction stream uses the Cont-In blocks to attend to part-aware context in the input image. Our intuition is that the pixel-aligned PAS images and the Cont-In block enables the network to determine the body part regions most salient for detecting emotion. Cont-In learns to modulate the network features by fusing the features of the intermediate layer with feature maps derived from PAS. Let $\mathbf{X} \in \mathbb{R}^{H \times W \times C}$ be the intermediate features from the $n - 1^{th}$ block of the base network. The PAS image \mathbf{P} is first passed through a series of convolutions and activation operations, denoted by $g(\cdot)$, to get an intermediate representation $\mathcal{G} = g(\mathbf{P})$ where $\mathcal{G} \in \mathbb{R}^{H \times W \times C}$. These feature maps are then concatenated with \mathbf{X} to get a fused representation $\mathbf{F} = \mathcal{G} \oplus \mathbf{X}$. \mathbf{F} is then passed through a second series of convolutions, activations, and finally batchnorm to get the feature map $\mathbf{X}' \in \mathbb{R}^{H \times W \times C'}$ which is then passed through to the n^{th} block of the base network (see Fig. 1).

4 Experiments

4.1 Experiment setup

Dataset and metrics. For the purpose of our experiments, we use the two-stream architecture of [22, 24] as our base implementation. We use their EMOTIC database [23], which is composed of images from MS-COCO [27], ADE20K [43] along with images downloaded from the web. The database offers two emotion representation labels; a set of 26 discrete emotional categories (Cat), and a set of three emotional dimensions, Valence, Arousal and Dominance from the VAD Emotional State Model [30]. Valence (V), is a measure of how positive or pleasant an emotion is (negative to positive); Arousal (A), is a measure of the agitation level of the person (non-active, calm, agitated, to ready to act); and Dominance (D) is a measure of the control level of the situation by the person (submissive, non-control, dominant, to in-control). The continuous dimensions (Cont) annotations of VAD are in a 1-10 scale.

Loss function. A dynamic weighted MSE loss L_{cat} is used on the Category classification output layer (Cat) of the model.

$$L_{cat} = \sum_{i=1}^{26} w_i (\hat{y}_i^{cat} - y_i^{cat})^2 \quad (4)$$

where i corresponds to 26 discrete categories shown in Tab. 1. \hat{y}_i^{cat} and y_i^{cat} are the prediction and ground-truth for the i^{th} category. The dynamic weight w_i are computed per batch and is based on the number of occurrences of each class in the batch. Since the occurrence of a particular class can be 0, [22, 24] defined an additional hyper-parameter c . The constant c is added to the dynamic weight w_i along with p_i , which is the probability of the i^{th} category. The final weight is defined as $w_i = \frac{1}{\ln(p_i + c)}$.

For the continuous (Cont) output layer, an L1 loss L_{cont} is employed.

$$L_{cont} = \frac{1}{C} \sum_{i=1}^C |\hat{y}_i^{cont} - y_i^{cont}| \quad (5)$$

Here i represents one of valence, arousal, and dominance (C). \hat{y}_i^{cont} and y_i^{cont} are the prediction and ground-truth for the i^{th} metric (VAD).

Baselines. We compare PERI to the SOTA baselines including *Emotic* Kosti et al. [22, 24], Huang et al. [19], Zhang et al. [42], Lei et al. [25], and Mittal et al. [32]. We reproduce the three stream architecture in [19] based on their proposed method. For a fair comparison, we compare PERI’s image-based model results with EmotiCon’s [32] image-based GCN implementation.

Implementation details. We use the two-stream architecture from Kosti et al. [22, 24]. Here both the image and body feature extraction streams are Resnet-18 [17] networks which are pre-trained on ImageNet [36]. All PAS images are re-sized to 128X128 similar to the input of the body feature extraction stream. The PAS image is created by plotting 501 landmarks $N + M$ on the base mask and passing it through a Gaussian filter of size $\sigma = 3$. We consider the same train, validation, and test splits provided by the EMOTIC [23] open repository.

Table 1: The average precision (AP) results on state-of-the-art methods and PERI. We see a consistent increase across a majority of the discrete class APs as well as the mAP using PERI.

Category	Kosti [22, 24]	Huang [19]	Zhang [42]	Lee [25]	Mittal [32]	PERI
Affection	28.06	26.45	46.89	19.90	36.78	38.87
Anger	6.22	6.52	10.87	11.50	14.92	16.47
Annoyance	12.06	13.31	11.23	16.40	18.45	20.61
Anticipation	93.55	93.31	62.64	53.05	68.12	94.70
Aversion	11.28	10.21	5.93	16.20	16.48	15.55
Confidence	74.19	74.47	72.49	32.34	59.23	78.92
Disapproval	13.32	12.84	11.28	16.04	21.21	21.48
Disconnection	30.07	30.07	26.91	22.80	25.17	36.64
Disquietment	16.41	15.12	16.94	17.19	16.41	18.46
Doubt/Confusion	15.62	14.44	18.68	28.98	33.15	20.36
Embarrassment	5.66	5.24	1.94	15.68	11.25	6.00
Engagement	96.68	96.41	88.56	46.58	90.45	97.93
Esteem	20.72	21.31	13.33	19.26	22.23	23.55
Excitement	72.04	71.42	71.89	35.26	82.21	79.21
Fatigue	7.51	8.74	13.26	13.04	19.15	13.94
Fear	5.82	5.76	4.21	10.41	11.32	7.86
Happiness	69.51	70.73	73.26	49.36	68.21	80.68
Pain	7.23	7.17	6.52	10.36	12.54	16.19
Peace	21.91	20.88	32.85	16.72	35.14	35.81
Pleasure	39.81	40.29	57.46	19.47	61.34	49.29
Sadness	7.60	8.04	25.42	11.45	26.15	18.32
Sensitivity	5.56	5.21	5.99	10.34	9.21	7.68
Suffering	6.26	7.83	23.39	11.68	22.81	19.85
Surprise	11.60	12.56	9.02	10.92	14.21	17.65
Sympathy	26.34	26.41	17.53	17.13	24.63	36.01
Yearning	10.83	10.86	10.55	9.79	12.23	15.32
mAP \uparrow	27.53	27.52	28.42	20.84	32.03	33.86

4.2 Quantitative results

Tab. 1 and Tab. 2 show quantitative comparisons between PERI and state-of-the-art approaches. Tab. 1 compares the average precision (AP) for each discrete

Table 2: The VAD and mean errors for continuous labels. The models include the state-of-the-art methods and PERI. We see a consistent decrease across each VAD $L1$ error along with the mean $L1$ error using PERI.

	Valence↓	Arousal↓	Dominance↓	VAD Error↓
Kosti et al. [22, 24]	0.71	0.91	0.89	0.84
Huang et al. [19]	0.72	0.90	0.88	0.83
Zhang et al. [42]	0.70	1.00	1.00	0.90
PERI	0.70	0.85	0.87	0.80

emotion category in the EMOTIC dataset [23]. Tab. 2 compares the valence, dominance and arousal $L1$ errors. Our model consistently outperforms existing approaches in both metrics. We achieve a significant 6.3% increase in mean AP (mAP) over our base network [22, 24] and a 1.8% improvement in mAP over the closest competing method [32]. Compared to methods that report VAD errors, PERI achieves lower mean and individual $L1$ errors and a 2.6% improvement in VAD error over our baseline [22, 24]. Thus, our results effectively shows that while only using pose or facial landmarks might lead to noisy gradients, especially in images with unreliable/occluded body or face, adding cues from both facial and body pose features where available lead to better emotional context. We further note that our proposed Cont-In Blocks are effective in reasoning about emotion context when comparing PERI with recent methods that use both body pose and facial landmarks [32].

4.3 Qualitative results

In order to understand the results further, we look at several visual examples, a subset of which are shown in Fig. 3. We choose Kosti et al. [22, 24] and Huang et al. [19] as our baselines as they are the closest SOTA methods.

We derive several key insights pertaining to our results. In comparison to Kosti et al. [22, 24] [19] and Huang et al., PERI fares better on examples where the face is clearly visible. This is expected as PERI specifically brings greater attention to facial features. Interestingly, our model also performs better for images where either the face or the body is partially visible (occluded/blurred). This supports our hypothesis that partial body poses as well as partial facial landmarks can supplement one another using our PAS image representation.

4.4 Ablation study

As shown in Tab. 3, we conduct a series of ablation experiments to create an optimal part-aware representation (PAS) and use the information in our base model effectively. For all experiments, we treat the implementation from Kosti et al. [22, 24] as our base network and build upon it.

PAS images. To get the best PAS representation, we vary the standard deviation (σ) of the Gaussian kernel applied to our PAS image. We show that

Table 3: Ablation studies. We divide this table into two sections. 1) Experiments to obtain the best PAS image. 2) Experiments to get the optimum method to use these PAS images (Cont-In blocks)

	mAP \uparrow	Valence \downarrow	Arousal \downarrow	Dominance \downarrow	Avg error \downarrow
Baselines					
Kosti et al. [22, 24]	27.53	71.16	90.95	88.63	83.58
Huang et al. [19]	27.52	72.22	89.92	88.31	83.48
PAS image experiments					
PAS $\sigma = 1$	33.32	70.77	83.75	89.47	81.33
PAS $\sigma = 3$	33.80	71.73	85.36	86.36	81.15
PAS $\sigma = 5$	33.46	70.36	87.56	85.39	81.10
PAS $\sigma = 7$	32.74	70.95	85.46	88.04	81.48
Cont-In block experiments					
Early fusion	32.96	70.60	85.95	87.59	81.38
Late fusion	32.35	71.73	85.60	87.12	81.48
Cont-In on both streams	29.30	72.43	87.23	89.97	83.21
PERI	33.86	70.77	84.56	87.49	80.94

$\sigma = 3$, gives us the best overall performance with a 5.9% increase in the mAP and a 2.5% decrease in the mean VAD error (Tab. 3: PAS image experiments) over the base network. From the use of PAS images, we see that retrieving context from input images that are aware of the facial landmarks and body poses is critical to achieving better emotion recognition performance from the base network.

Experimenting with Cont-In blocks. To show the effectiveness of Cont-In blocks, we compare its performance with early and late fusion in Tab. 3. For early fusion, we concatenate the PAS image as an additional channel to the body-crop image in the body feature extraction stream. For late fusion, we concatenate the fused output of the body and image feature extraction streams with the downsampled PAS image. As opposed to PERI, we see a decline in performance for both mAP and VAD error when considering early and late fusion. From this we conclude that context infusion at intermediate blocks is important for accurate emotion recognition.

Additionally, we considered concatenating the PAS images directly to the intermediate features instead of using a Cont-In block. However, feature concatenation in the intermediate layers changes the backbone ResNet architecture, severely limiting gains from ImageNet [36] pretraining. This is apparent in the decrease in performance from early fusion, which may be explained, in part, by the inability to load ImageNet weights in the input layer of the backbone network. In contrast, Cont-In blocks are fully compatible with any emotion recognition network and do not alter the network backbone.

In the final experiment, we added Cont-In blocks to both the image feature extraction stream and the body feature extraction stream. Here we discovered that if we regulate the intermediate features of both streams as opposed to just the body stream the performance declines. A possible reason could be that

contextual information from a single person does generalise well to the entire image with multiple people.

PERI. From our ablation experiments, we found that PERI works best overall. It has the highest mAP among the ablation experiments as well as a lowest mean $L1$ error for VAD. While there are other hyper-parameters that have better $L1$ errors for Valence, Arousal, and Dominance independently, (different Gaussian standard deviations (σ_k)), these hyper-parameters tend to perform worse overall compared to PERI.

5 Conclusion

Existing methods for in the wild emotion recognition primarily focus on either face or body, resulting in failures under challenging scenarios such as low resolution, occlusion, blur etc. To address these issues, we introduce PERI, a method that effectively represents body poses and facial landmarks together in a pixel-aligned part aware contextual image representation, PAS. We argue that using both features results in complementary information which is effective in challenging scenarios. Consequently, we show that PAS allows better emotion recognition not just in examples with fully visible face and body features, but also when one of the two features sets are missing, unreliable or partially available.

To seamlessly integrate the PAS images with a baseline emotion recognition network, we introduce a novel method for modulating intermediate features with the part-aware spatial (PAS) context by using context infusion (Cont-In) blocks. We demonstrate that using Cont-In blocks works better than a simple early or late fusion. PERI significantly outperforms the baseline emotion recognition network of Kosti et al. [23, 24]. PERI also improves upon existing state-of-the-art methods on both the mAP and VAD error metrics.

While our method is robust towards multiple in-the-wild challenging scenarios, we do not model multiple-human scenes and dynamic environments. In the future, we wish to further extend Cont-In blocks to utilise the PAS context better. Instead of modeling explicit attention using PAS images, it might be interesting to learn part-attention implicitly using self and cross-attention blocks, but we leave this for future work. Additionally, we also seek to explore multi-modal input beyond images, such as depth, text and audio.

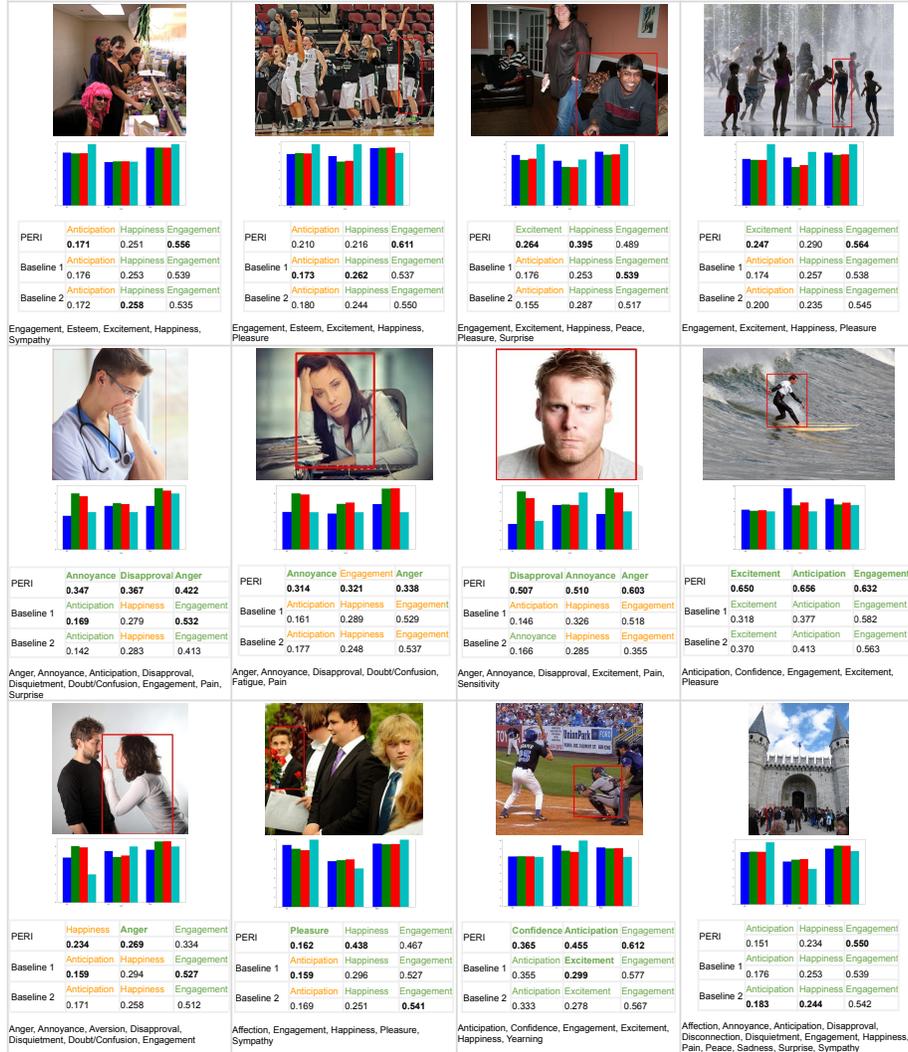


Fig. 3: The figure shows visual results on selected examples. Here each example is separated into 3 parts: the original image with the person of interest in a bounding box; the VAD bar plot (scale 1-10); and the top 3 emotion categories predicted by each model. For VAD values the blue, green, red, cyan columns correspond to PERI, baseline 1 (Kosti et al. [22, 24]), baseline 2 (Huang et al. [19]) and the ground-truth value respectively. For the predicted categories we highlight the category in green if they are present in the ground-truth and orange if they aren't. The ground-truth categories associated with each example are written as a list below the predictions.

References

1. Ahmed, F., Bari, A.S.M.H., Gavrilova, M.L.: Emotion recognition from body movement. *IEEE Access* **8**, 11761–11781 (2020). <https://doi.org/10.1109/ACCESS.2019.2963113> **2**
2. Benitez-Quiroz, C.F., Srinivasan, R., Martinez, A.M.: Emotionet: An accurate, real-time algorithm for the automatic annotation of a million facial expressions in the wild. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 5562–5570 (2016). <https://doi.org/10.1109/CVPR.2016.600> **2, 4**
3. Castillo, G., Neff, M.: What do we express without knowing?: Emotion in gesture. In: AAMAS (2019) **4**
4. Chu, W.S., De la Torre, F., Cohn, J.F.: Selective transfer machine for personalized facial expression analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **39**(3), 529–545 (2017). <https://doi.org/10.1109/TPAMI.2016.2547397> **4**
5. Coulson, M.: Attributing emotion to static body postures: Recognition accuracy, confusions, and viewpoint dependence. *Journal of nonverbal behavior* **28**(2), 117–139 (2004) **2**
6. Crenn, A., Khan, R.A., Meyer, A., Bouakaz, S.: Body expression recognition from animated 3d skeleton. In: 2016 International Conference on 3D Imaging (IC3D). pp. 1–7 (2016). <https://doi.org/10.1109/IC3D.2016.7823448> **4**
7. De Gelder, B., Van den Stock, J.: The bodily expressive action stimulus test (beast). construction and validation of a stimulus basis for measuring perception of whole body expression of emotions. *Frontiers in psychology* **2**, 181 (2011) **2**
8. Duncan, D., Shine, G., English, C.: Facial emotion recognition in real time. *Computer Science* pp. 1–7 (2016) **2**
9. Ekman, P., Friesen, W.V.: Constants across cultures in the face and emotion. *Journal of personality and social psychology* **17**(2), 124 (1971) **4**
10. Ekman, P., Friesen, W.V.: Facial action coding system: a technique for the measurement of facial movement (1978) **4**
11. Eleftheriadis, S., Rudovic, O., Pantic, M.: Discriminative shared gaussian processes for multiview and view-invariant facial expression recognition. *IEEE Transactions on Image Processing* **24**(1), 189–204 (2015). <https://doi.org/10.1109/TIP.2014.2375634> **4**
12. Eleftheriadis, S., Rudovic, O., Pantic, M.: Joint facial action unit detection and feature fusion: A multi-conditional learning approach. *IEEE Transactions on Image Processing* **25**(12), 5727–5742 (2016). <https://doi.org/10.1109/TIP.2016.2615288> **2, 4**
13. de Gelder, B.: Towards the neurobiology of emotional body language. *Nature Reviews Neuroscience* **7**, 242–249 (2006) **4**
14. Gross, M.M., Crane, E.A., Fredrickson, B.L.: Effort-shape and kinematic assessment of bodily expression of emotion during gait. *Human movement science* **31** **1**, 202–21 (2012) **4**
15. Gunes, H., Piccardi, M.: Bi-modal emotion recognition from expressive face and body gestures. *J. Netw. Comput. Appl.* **30**, 1334–1345 (2007) **3, 4**
16. Gupta, A., Agrawal, D., Chauhan, H., Dolz, J., Pedersoli, M.: An attention model for group-level emotion recognition. In: Proceedings of the 20th ACM International Conference on Multimodal Interaction. pp. 611–615 (2018) **5**
17. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. *CoRR abs/1512.03385* (2015), <http://arxiv.org/abs/1512.03385> **7, 10**

18. Hu, M., Wang, H., Wang, X., Yang, J., Wang, R.: Video facial emotion recognition based on local enhanced motion history image and cnn-ctslstm networks. *Journal of Visual Communication and Image Representation* **59**, 176–185 (2019) [2](#)
19. Huang, Y., Wen, H., Qing, L., Jin, R., Xiao, L.: Emotion recognition based on body and context fusion in the wild. In: 2021 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW). pp. 3602–3610 (2021). <https://doi.org/10.1109/ICCVW54120.2021.00403> [2](#), [4](#), [9](#), [10](#), [11](#), [12](#), [14](#)
20. Kleinsmith, A., Bianchi-Berthouze, N.: Recognizing affective dimensions from body posture. In: *International conference on affective computing and intelligent interaction*. pp. 48–58. Springer (2007) [2](#)
21. Ko, B.C.: A brief review of facial emotion recognition based on visual information. *sensors* **18**(2), 401 (2018) [2](#)
22. Kosti, R., Alvarez, J.M., Recasens, A., Lapedriza, A.: Context based emotion recognition using emotic dataset. arXiv preprint arXiv:2003.13401 (2020) [2](#), [4](#), [5](#), [7](#), [9](#), [10](#), [11](#), [12](#), [14](#)
23. Kosti, R., Álvarez, J.M., Recasens, A., Lapedriza, À.: Emotic: Emotions in context dataset. 2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW) pp. 2309–2317 (2017) [3](#), [4](#), [5](#), [9](#), [10](#), [11](#), [13](#)
24. Kosti, R., Álvarez, J.M., Recasens, A., Lapedriza, À.: Emotion recognition in context. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) pp. 1960–1968 (2017) [2](#), [4](#), [5](#), [7](#), [9](#), [10](#), [11](#), [12](#), [13](#), [14](#)
25. Lee, J., Kim, S., Kim, S., Park, J., Sohn, K.: Context-aware emotion recognition networks. In: 2019 IEEE/CVF International Conference on Computer Vision (ICCV). pp. 10142–10151 (2019). <https://doi.org/10.1109/ICCV.2019.01024> [5](#), [9](#), [10](#)
26. Li, Z., Imai, J.i., Kaneko, M.: Facial-component-based bag of words and phog descriptor for facial expression recognition. In: 2009 IEEE International Conference on Systems, Man and Cybernetics. pp. 1353–1358 (2009). <https://doi.org/10.1109/ICSMC.2009.5346254> [4](#)
27. Lin, T., Maire, M., Belongie, S.J., Bourdev, L.D., Girshick, R.B., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft COCO: common objects in context. *CoRR abs/1405.0312* (2014), <http://arxiv.org/abs/1405.0312> [9](#)
28. Lugaresi, C., Tang, J., Nash, H., McClanahan, C., Uboweja, E., Hays, M., Zhang, F., Chang, C., Yong, M.G., Lee, J., Chang, W., Hua, W., Georg, M., Grundmann, M.: Mediapipe: A framework for building perception pipelines. *CoRR abs/1906.08172* (2019), <http://arxiv.org/abs/1906.08172> [5](#), [6](#)
29. Mehendale, N.: Facial emotion recognition using convolutional neural networks (ferc). *SN Applied Sciences* **2**(3), 1–8 (2020) [2](#)
30. Mehrabian, A.: Framework for a comprehensive description and measurement of emotional states. *Genetic, social, and general psychology monographs* **121** **3**, 339–61 (1995) [4](#), [7](#), [9](#)
31. Mellouk, W., Handouzi, W.: Facial emotion recognition using deep learning: review and insights. *Procedia Computer Science* **175**, 689–694 (2020) [2](#)
32. Mittal, T., Guhan, P., Bhattacharya, U., Chandra, R., Bera, A., Manocha, D.: Emoticon: Context-aware multimodal emotion recognition using frege’s principle. 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 14222–14231 (2020) [4](#), [5](#), [9](#), [10](#), [11](#)
33. Öhman, A., Dimberg, U.: Facial expressions as conditioned stimuli for electrodermal responses: a case of ‘preparedness’? *Journal of Personality and Social Psychology* **36**(11), 1251 (1978) [2](#), [4](#)

34. Pantic, M., Rothkrantz, L.: Expert system for automatic analysis of facial expression. *Image and Vision Computing* **18**, 881–905 (08 2000). [https://doi.org/10.1016/S0262-8856\(00\)00034-2](https://doi.org/10.1016/S0262-8856(00)00034-2) 4
35. Pranav, E., Kamal, S., Chandran, C.S., Supriya, M.: Facial emotion recognition using deep convolutional neural network. In: 2020 6th International conference on advanced computing and communication Systems (ICACCS). pp. 317–320. IEEE (2020) 2
36. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)* **115**(3), 211–252 (2015). <https://doi.org/10.1007/s11263-015-0816-y> 10, 12
37. Russell, J.A., Bullock, M.: Multidimensional scaling of emotional facial expressions: similarity from preschoolers to adults. *Journal of personality and social psychology* **48**(5), 1290 (1985) 2, 4
38. Shen, Z., Cheng, J., Hu, X., Dong, Q.: Emotion recognition based on multi-view body gestures. In: 2019 IEEE International Conference on Image Processing (ICIP). pp. 3317–3321 (2019). <https://doi.org/10.1109/ICIP.2019.8803460> 2
39. Tümen, V., Söylemez, Ö.F., Ergen, B.: Facial emotion recognition on a dataset using convolutional neural network. In: 2017 International Artificial Intelligence and Data Processing Symposium (IDAP). pp. 1–5. IEEE (2017) 2
40. Wu, J., Zhang, Y., Zhao, X., Gao, W.: A generalized zero-shot framework for emotion recognition from body gestures (2020). <https://doi.org/10.48550/ARXIV.2010.06362>, <https://arxiv.org/abs/2010.06362> 4, 5
41. Zacharatos, H., Gatzoulis, C., Chrysanthou, Y.L.: Automatic emotion recognition based on body movement analysis: a survey. *IEEE computer graphics and applications* **34**(6), 35–45 (2014) 2
42. Zhang, M., Liang, Y., Ma, H.: Context-aware affective graph reasoning for emotion recognition. 2019 IEEE International Conference on Multimedia and Expo (ICME) pp. 151–156 (2019) 4, 5, 9, 10, 11
43. Zhou, B., Zhao, H., Puig, X., Fidler, S., Barriuso, A., Torralba, A.: Scene parsing through ade20k dataset. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2017) 9