# Detecting Emotion from People and Scenery: Midterm Update

Logan Preston (lpreston2@wisc.edu), Neal Desai (nbdesai2@wisc.edu),
and Nikhil Nigam (nnigam@wisc.edu)

## 1    Introduction

The motivation for this project is to investigate approaches for emotion detection and recognition. This builds on existing work of facial recognition and tracking to predict what the user is likely feeling based based on image-specific context. Existing work often looks to face expressions, body language, and scene context. These factors combined with individual variation in how people express emotion makes this a challenging but useful problem to solve.

Specifically, there is social value in having a good solution to the emotion recognition problem. The use of AI-infused systems is growing in daily life and generally is expected to continue that trend. Human AI interaction guidelines state that AI infused systems should have "socially appropriate behaviors" [5] or "match relevant social norms" [1], showing that being aware of social expectations is helpful. Detecting user emotions can inform the AI system on appropriate social expectations and also potentially detect emotional anomalies. This could be used to make artificial intelligence applications more socially aware, friendly, or provide better support to humans. For example, identifying if a user is frustrated may prompt the system to check if the user has a question, provide supportive reminders, or reduce the amount of notifications it sends to the user. In general, we expect identifying emotions accurately will enhance interactions between AI infused systems and their users.

## 2    State of the Art

This problem is especially thought-provoking because we already know computers can identify faces, track object paths, and other related tasks. However, the work on emotion detection is not as well-established. Classical emotion recognition techniques take a more localized approach by simply focusing on the prominent facial features. These techniques have been now bettered, if not matched, by deep learning algorithms such as CNNs or LSTMs [6].

Thus, prior work reviews emotions by looking at the faces alone rather than considering body language and scene context. Some work adds in context from the surrounding the scene [8] but efforts in this area seem minimal and don't focus on quantifying the potential improvement of considering scene context vs just faces. We are interested in seeing how important the overall scene context is for the accuracy of the model or if the face / body language provides the vast majority of the data needed for accurate classification. Apart from the further exploration by those who initially included the scene for analysis [7], another method adopted by researchers involved masking the face from the image and treating facial features and the scene separately and adopting ensemble learning techniques for their prediction [9]. These findings act as a starting point for a more nuanced technique which can be developed and researched further.

## 3    Proposed Solution

Our solution will start by identifying people from the image, which will require segmenting the image as a prerequisite. After the humans in each image are identified, we can look to them for potential clues into the emotions using their faces and/or body language. We will further investigate the impact of incorporating the full body and scene context for our emotion recognition as opposed to just isolated faces.

We will evaluate the performance of solutions using only expressions and body language and compare those results to the performance of solutions that also include information on the scene. This will identify the impact of the scene information and determine how valuable of information it is to analyze compared to the individuals in a scene.

To train our model, we plan to use the Emotic dataset. The dataset contains images of several people in real environments. Each person is appropriately annotated with their emotions. The dataset offers both a discrete analysis using 26 distinct emotion categories as well as

a continuous dataset that seeks to quantify each emotion into three dimensions, namely valence, arousal, and dominance [7]. We will incorporate both the discrete and continuous cases in our analysis. We plan on initially using Convolutional Neural Networks to accomplish this classification and will explore other deep learning model architectures that may improve upon the baseline CNN approach.

# 4 Progress and Difficulties

## 4.1 Progress

We have reviewed the Emotic data set, created utilities to parse it, and have used it in an implementation of the Baseline model that can infer the emotion in an image. The model can separate out scene from people, and we have trained it on a **small portion** of the data set.

The baseline model consists of 3 primary components which drive the classification. The first component takes the region of the image comprising the person whose feelings are to be estimated and extracts its most relevant features. The second component takes as input the entire image and extracts global features for providing the necessary contextual support. Finally, the third component is a fusion network that takes as input the image and body features and estimates the discrete categories and the continuous dimensions [8]. A high level overview of the architecture can be seen in Figure 1 for the body and scene split.
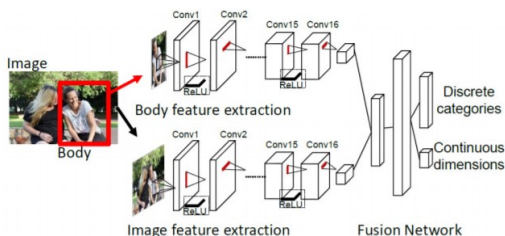


Figure 1: Overall Architecture for Body/Scene Split

Given limited computational resources, we have resorted to using the popular Resnet-18 architecture for the first two components [4]. The feature extractors are initialized using pre-trained models from large scale datasets such as ImageNet [2] and Places [10]. ImageNet consists of a range of objects (including people), so incorporating that will improve our ability to learn the region of significance for a target person in the image. Places is a data set specifically created for high level visual understanding tasks such as recognizing scene categories. Pre-training the image feature extraction

model using this data set ensures providing global (high level) contextual support.

For our final report, we expect to have a contribution amount for the scene, the body, and the face based on the weights for each component. Calculating the contribution by basing it on the weights is a naive approach, but we expect to see some trend, and may update this in the final proposal if we identify a less naive approach. We will also consider metrics as part of our train / test split which quantify the classification error (in the discrete case) or mean-squared-error (in the continuous case) to gain an understanding of how the actual task would be impacted by the difference in weights.

We also have created our project website to include much of the information here, it can be found on Github[3].

## 4.2 Example Results

Figure 2 is the training loss for the original model and Figure 3 is our trained model, showing that our training closely matches the original implementation.
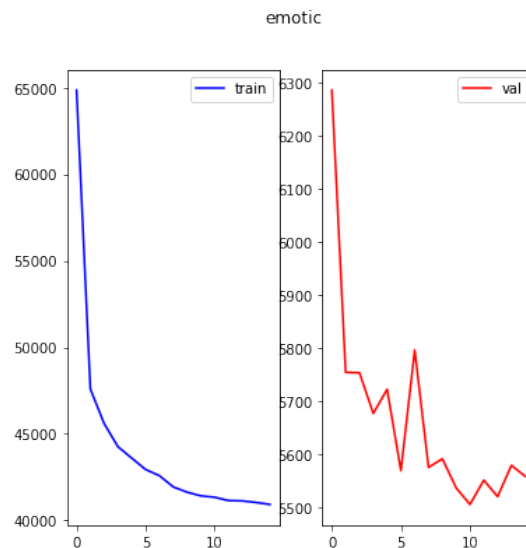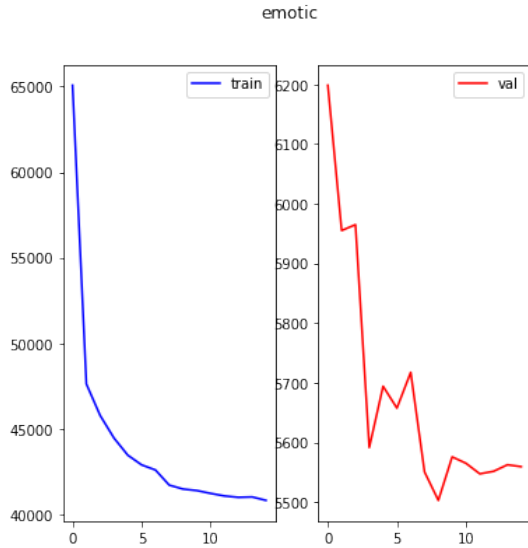


Figure 2: Training Loss vs Time, Original

Figure 3: Training Loss vs Time, Our Implementation

## 4.3 Difficulties

Parsing our data set did prove a little difficult due to some weak documentation around it; however, we found a number of helpful data elements such as the coordinates of the box that surround the portion of the image that's being evaluated. An example of an image can be seen in Figure 4, and with this bounding box drawn on in Figure 5. This way we can identify which person is matched to each emotion in the data set for the model.

This bounding box annotations mean we don't implement a custom computer vision algorithm to separate out the individual in the image at this point. We only feed the separate parts of the image to the machine learning algorithms. We expect to add more discrete computer vision functionality to identify the locations of the faces, such as utilizing Haar feature detection, for the final report. We will account for these bounding boxes that are given by the data set to ensure we have the correct face in the image to match with the emotion from the data set.

## 4.4 Changes to Proposal

We have not had any major changes to our proposal at this time. We have confirmed our initial ideas of how to measure the impact of each of the image components discussed previously. To quantify this, we will use a combination of the weights as well as pertinent error metrics to calculate the contribution for the scene, body, and face individually.



Figure 4: A sample image from the Emotic Dataset



Figure 5: A sample image from the Emotic Dataset, with a bounding box drawn to identify the person

## 5 Updated Time Table

Table 1 is an updated time table for our project, unchanged from the first report and trimmed down to remaining tasks.

| Objective | Due Date |
|---|---|
| Extend testing with complete data set to identify impact of scene vs body language vs face | 4/15 |
| Begin creating final report and presentation | 4/15 |
| Finalize final report and presentation | 4/25 |

Table 1: Current Time Table

## References

[1] S. Amershi, D. Weld, M. Vorvoreanu, A. Fourney, B. Nushi, P. Collisson, J. Suh, S. Iqbal, P. Bennett,

K. Inkpen, J. Teevan, R. Kikin-Gil, and E. Horvitz. Guidelines for human-ai interaction. In *CHI 2019*. ACM, May 2019. CHI 2019 Honorable Mention Award.

[2] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009.

[3] N. Desai, N. Nigam, and L. Preston. Comp sci 766 term project page. https://loganpreston.github.io, 2022.

[4] 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)*, pages 770–778, 2016.

[5] E. Horvitz. Principles of mixed-initiative user interfaces. In *Proceedings of CHI '99, ACM SIGCHI Conference on Human Factors in Computing Systems, Pittsburgh, PA, ACM Press.*, pages 159–166, May 1999.

[6] B. C. Ko. A brief review of facial emotion recognition based on visual information. volume 18, 2018.

[7] R. Kosti, J. Alvarez, A. Recasens, and A. Lapedriza. Context based emotion recognition using emotic dataset. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019.

[8] R. Kosti, J. M. Alvarez, A. Recasens, and A. Lapedriza. Emotion recognition in context. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

[9] J. Lee, S. Kim, S. Kim, J. Park, and K. Sohn. Context-aware emotion recognition networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.

[10] B. Zhou, A. Khosla, À. Lapedriza, A. Torralba, and A. Oliva. Places: An image database for deep scene understanding. volume abs/1610.02055, 2016.