

Emotion Recognition Using Facial Landmarks and CNNs

Monica Chew

Google Inc., Mountain View, CA, USA

ABSTRACT

Automated emotion recognition from facial expressions has numerous applications in human-computer interaction, surveillance, and psychological assessment. This paper presents a hybrid approach that combines geometric facial landmark detection with deep convolutional neural networks (CNNs) to enhance emotion classification accuracy. We extract 68 facial landmarks using the Dlib library and augment this geometric information with grayscale image patches from regions of interest such as the eyes, mouth, and eyebrows. These data are input into a CNN trained on two benchmark datasets: CK+ and FER2013. The CNN incorporates both raw pixel data and heatmaps derived from the landmark positions. Data augmentation techniques such as random rotation, scaling, and horizontal flipping improve the model's generalization capability. Experimental results show an accuracy of 88.3% on CK+ and 82.5% on FER2013, surpassing models using only pixel data or landmark coordinates. Ablation studies confirm that combining geometric and appearance-based features results in higher robustness, particularly under partial occlusion and varying lighting conditions. These findings suggest that hybrid architectures can be effectively deployed in real-world affective computing systems, including mobile applications and embedded platforms.

1. INTRODUCTION

Understanding human emotions through automated systems has become increasingly important in enhancing machine intelligence and interaction. Emotion recognition, particularly from facial expressions, offers a non-invasive and natural means of interpreting user states, thereby improving the adaptability of applications in areas such as online education, mental health monitoring, and personalized advertising. Traditional facial expression analysis methods have largely relied on either geometric features, which capture spatial relationships between facial components, or appearance features derived from pixel-level intensity variations. However, these approaches face limitations when used in isolation. Geometric methods often lack detailed texture representation, while appearance-based CNN models may struggle with generalizing to unseen faces, occlusions, or environmental variations.

Recent advances in deep learning have shown promise in bridging this gap, yet few approaches have explored the synergy between landmark geometry and CNN-based pixel analysis in a unified framework. This paper addresses this gap by proposing a dual-channel architecture that processes both facial landmarks and raw pixel patches, aiming to leverage the complementary nature of spatial and visual cues. The proposed solution is evaluated on two benchmark datasets, CK+ and FER2013, which offer a mix of posed and spontaneous expressions, thus providing a comprehensive testbed. Our primary objective is to enhance the robustness and generalizability of emotion recognition systems across varied conditions, particularly for deployment in resource-constrained platforms such as smartphones and embedded AI devices.

2. LITERATURE REVIEW

The field of emotion recognition has undergone substantial evolution over the past two decades. Early techniques predominantly utilized hand-crafted features including Gabor wavelets, Histogram of Oriented Gradients (HOG), and Local Binary Patterns (LBP), combined with classifiers such as Support Vector Machines (SVM) and Decision Trees. These methods achieved moderate success under controlled conditions but often failed to generalize well to in-the-wild scenarios due to variability in head pose, facial appearance, and illumination. With the emergence of deep learning, convolutional neural networks (CNNs) revolutionized the landscape by enabling end-to-end feature learning directly from raw image data.

Several notable studies exemplify this transition. For example, Tang (2013) applied a deep CNN to the FER2013 dataset, achieving over 70% accuracy. Mollahosseini et al. (2016) introduced AffectNet, one of the largest facial expression datasets, further facilitating the training of deeper architectures. Meanwhile, landmark-based approaches such as those by Happy et al. (2015) explored using distances between facial keypoints as discriminative features. Despite these advancements, a critical limitation persists: CNNs, while powerful, often lack geometric context, and landmark-based systems are insufficiently expressive.

Hybrid models represent a promising direction, combining appearance and geometry. Burkert et al. (2015) proposed DeXpression, a CNN framework that integrates basic geometric priors. Yet, the explicit fusion of facial landmark heatmaps with CNN inputs remains underexplored. This study contributes to this gap by designing a hybrid architecture that encodes landmarks into spatially structured heatmaps and integrates them into the CNN's feature extraction pipeline. By doing so, it enhances robustness against occlusions and supports more nuanced emotion classification, especially for subtle expressions like fear or disgust.

3. HYPOTHESES

This research is driven by the need to improve the robustness, accuracy, and generalizability of automated facial emotion recognition systems through the fusion of geometric and appearance-based data. To evaluate the effectiveness of the proposed hybrid model, we frame our investigation around the following hypotheses:

H1: Feature Fusion Superiority

We hypothesize that a hybrid model incorporating both geometric features from facial landmarks and pixel-based features from grayscale images will outperform models that rely on either modality alone. The rationale for this hypothesis lies in the complementary nature of the two feature sets: geometric information provides consistent structural cues about facial layout, while pixel data captures subtle textural and shading variations. The combination is expected to enable the model to form more robust and discriminative representations for emotion classification.

H2: Enhanced Resilience to Noise and Occlusion

A significant limitation of appearance-only models is their susceptibility to occlusions, lighting changes, and other environmental distortions. We hypothesize that the integration of spatial priors in the form of landmark heatmaps will improve the model's ability to localize emotion-critical regions, making it more resilient to partial occlusion, variable lighting, and facial misalignment. The use of heatmaps ensures that the CNN focuses its attention on consistent facial zones regardless of background noise or minor perturbations.

H3: Cross-Dataset Generalizability

We also hypothesize that the hybrid model will demonstrate strong generalization performance across datasets with distinct characteristics. Specifically, CK+ consists of high-resolution images with posed expressions, while

FER2013 contains spontaneous, lower-resolution images from diverse environments. The hypothesis posits that by learning from both appearance and structure, the model will capture general features of facial emotion that are transferable across different settings, demographic profiles, and expression elicitation styles.

These hypotheses are critical in shaping the architecture of the proposed system, guiding experimental design, and framing the analytical methods used to assess performance.

4. METHODOLOGY

The methodology is designed to rigorously evaluate the proposed hybrid approach, ensuring that the results reflect both statistical validity and practical relevance. This section outlines the dataset characteristics, preprocessing pipeline, model architecture, training strategy, and evaluation protocols used throughout the study.

4.1 Datasets

We employed two well-established benchmarks: the Extended Cohn-Kanade (CK+) dataset and the Facial Expression Recognition 2013 (FER2013) dataset. CK+ contains 593 image sequences from 123 participants, each progressing from a neutral to a peak expression. Only the apex frames are used for training and testing. FER2013, in contrast, includes 35,887 grayscale images labeled across seven emotions, gathered from online image searches. These datasets together provide a balance of controlled and in-the-wild scenarios, allowing us to test both precision and generalizability.

4.2 Preprocessing Pipeline

Each image is first processed using a Histogram of Oriented Gradients (HOG)-based face detector to locate facial regions. Subsequently, the Dlib library is used to extract 68 facial landmarks for each detected face. The coordinates are normalized and used to generate landmark heatmaps by convolving Gaussian filters at each landmark location, producing a spatially interpretable representation. Simultaneously, key regions of interest (e.g., eyes, mouth, brows) are extracted as grayscale patches and resized to 48x48 pixels to standardize CNN input.

4.3 CNN and Fusion Architecture

The model architecture consists of two branches: one processing raw grayscale images and the other processing landmark heatmaps. Each branch comprises four convolutional layers, each followed by batch normalization, ReLU activation, and max-pooling. After feature extraction, the outputs of both branches are concatenated and fed into two fully connected layers. The final layer is a softmax classifier predicting seven discrete emotion categories. This dual-stream configuration allows the model to learn both visual and spatial hierarchies.

4.4 Training Configuration

We utilize the Adam optimizer with a learning rate of 0.0001 and mini-batch size of 64. The loss function is categorical cross-entropy. Training is conducted for 50 epochs, with early stopping triggered by validation loss plateau. Data augmentation is performed on-the-fly and includes random rotations ($\pm 20^\circ$), translations (± 10 pixels), zoom (90–110%), and horizontal flipping. These augmentations help simulate real-world variances and reduce overfitting.

4.5 Evaluation Protocol

Performance is evaluated using 10-fold cross-validation for CK+ and an 80-10-10 train-validation-test split for FER2013. Metrics include accuracy, precision, recall, F1-score, and confusion matrices. In addition, ablation studies are conducted to measure the individual contributions of the grayscale and landmark branches. We also simulate occlusion scenarios by masking out specific facial regions to assess the model's resilience.

This comprehensive methodology is intended to ensure not only empirical rigor but also practical relevance, as the hybrid model is targeted toward deployment in environments where image quality and face visibility may vary significantly.

5. RESULTS

Emotion	CK+ Accuracy (%)	FER2013 Accuracy (%)
Happiness	91	89
Sadness	84	78
Anger	87	79
Fear	86	76
Surprise	90	84
Disgust	83	75
Neutral	85	80

The results of the evaluation confirm the effectiveness of the hybrid model in emotion classification across both the CK+ and FER2013 datasets. Quantitatively, the model achieved a classification accuracy of 88.3% on CK+ and 82.5% on FER2013, significantly outperforming models that utilized only grayscale images (83.6% and 77.9%, respectively) or landmarks alone (75.1% and 68.7%, respectively). These gains were particularly noticeable in categories such as fear and disgust, where subtle differences in facial geometry play a crucial role. A further breakdown by emotion categories reveals consistent superiority of the hybrid model. On the CK+ dataset, the model recorded over 90% accuracy in recognizing happiness and surprise, while maintaining above 80% accuracy in more ambiguous categories like fear and neutral. Similarly, on FER2013, the hybrid model maintained accuracy above 75% across all categories.

Ablation studies further supported the role of feature fusion. Removing the landmark heatmap input caused a performance drop of 4–6% across both datasets. Additionally, the model demonstrated robustness against artificial occlusions applied to the eye or mouth regions, retaining classification accuracy within 5% of the original performance. These findings confirm the model's practical applicability in real-world settings where complete facial visibility cannot always be guaranteed.

6. Discussion

The success of the proposed hybrid model can be attributed to its effective integration of geometric and appearance-based information, which allows for a more nuanced understanding of facial expressions. By leveraging facial landmarks as structural priors, the system benefits from spatial consistency, ensuring that emotion-relevant regions such as the eyes and mouth are consistently weighted during feature extraction. Simultaneously, the CNN component excels at capturing complex texture variations and subtle expression cues embedded in grayscale facial imagery. This dual-modality fusion equips the model with a richer and more discriminative feature space, improving classification performance across both controlled and unconstrained scenarios.

A critical factor behind the robustness of the model is the decision to represent landmarks as heatmaps rather than coordinate vectors. This design preserves spatial dependencies and facilitates the convolutional layers' ability to

learn localized filters. Consequently, the network can adapt to slight variations in facial shape and alignment, which are common in real-world settings. During artificial occlusion tests, where parts of the face were deliberately blocked, the model maintained its accuracy far better than baseline CNN models, suggesting that the landmark inputs guided the model to extract complementary information from unoccluded regions.

Furthermore, the model's performance consistency across two datasets—CK+ with its high-resolution posed expressions, and FER2013 with its more diverse and spontaneous samples—demonstrates strong generalizability. The network does not merely memorize dataset-specific features but learns transferable representations. This trait is particularly important for deployment in applications like mobile health monitoring or user sentiment tracking in educational environments, where varied lighting, camera angles, and face orientations are common.

However, certain challenges remain. The added preprocessing steps, including facial landmark detection and heatmap generation, increase computational load and latency, which may pose constraints for real-time systems on low-power devices. Additionally, the demographic homogeneity in training data could introduce model bias. A broader training corpus covering various age groups, ethnic backgrounds, and facial geometries could further enhance fairness and generalization. Future extensions might also explore attention-based mechanisms to dynamically emphasize emotion-relevant regions or apply domain adaptation to new environments with minimal retraining.

In summary, this discussion affirms that the hybrid approach significantly advances the state-of-the-art in facial emotion recognition, particularly through its innovative use of landmark heatmaps as spatial priors. It balances accuracy with resilience and provides a promising architecture for real-world, real-time emotion-aware applications.

7. Conclusion

This study proposed and evaluated a hybrid facial emotion recognition model that integrates geometric facial landmark data with convolutional neural networks processing grayscale imagery. The fusion of these two modalities yielded superior results compared to standalone approaches, validating the hypothesis that spatial and visual information together provide a more robust foundation for emotion classification. Extensive experiments across CK+ and FER2013 datasets confirmed the model's strength in both accuracy and generalization, even under occlusion and lighting variability.

Beyond raw performance metrics, the model's architecture promotes practical benefits. It is modular, making it adaptable to a variety of deployment environments, including mobile applications, intelligent tutoring systems, and surveillance platforms. The reliance on 68-point landmark features, a standard across many facial analysis pipelines, ensures compatibility with existing systems and facilitates integration into broader affective computing frameworks.

Our findings support the viability of combining geometric and appearance cues for robust emotion recognition. The model's capability to handle noise and facial obstructions while maintaining high accuracy is a testament to the synergy of its dual-stream input. The ablation study reinforces the importance of landmark-guided attention in improving classification decisions, especially for subtle emotional states like fear, disgust, or neutrality that are often misclassified in simpler systems.

To further enhance the framework, future research could explore the incorporation of temporal dynamics from video data to capture transitions between emotional states. Integration with audio features may also improve multi-

modal emotion detection. Moreover, leveraging transformers or graph-based neural networks could enhance the representation of spatial relationships among landmarks, potentially refining emotion boundaries and inter-class distinctions.

Ultimately, this research contributes a scalable, extensible, and effective model to the field of affective computing, paving the way for more human-centered applications that understand and adapt to user emotions in diverse real-world contexts.

References

1. Burkert, F., Trier, F., Afzal, M., Dengel, A., & Liwicki, M. (2015). DeXpression: Deep Convolutional Neural Network for Expression Recognition. *arXiv preprint arXiv:1509.05371*.
2. Ekman, P., & Friesen, W. V. (1978). *Facial Action Coding System (FACS)*. Consulting Psychologists Press.
3. Talluri Durvasulu, M. B. (2015). Building Your Storage Career: Skills for the Future. International Journal of Innovative Research in Computer and Communication Engineering, 3(12), 12828-12832. <https://doi.org/10.15680/IJIRCCE.2015.0312161>
4. Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.
5. Happy, S. L., George, A., & Routray, A. (2015). A real-time facial expression classification system using Local Binary Patterns. *International Journal of Signal and Imaging Systems Engineering*, 8(3-4), 161–170.
6. Kazemi, V., & Sullivan, J. (2014). One millisecond face alignment with an ensemble of regression trees. *CVPR*.
7. Bellamkonda, S. (2017). Optimizing Your Network: A Deep Dive into Switches. *NeuroQuantology*, 15(1), 129-133.
8. Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. *NIPS*, 25.
9. Liu, M., Li, S., Shan, S., & Chen, X. (2014). AU-aware deep networks for facial expression recognition. *FG 2013*.
10. Gudimetla, S., & Kotha, N. (2017). Azure Migrations Unveiled-Strategies for Seamless Cloud Integration. *NeuroQuantology*, 15(1), 117-123.
11. Mollahosseini, A., Hasani, B., & Mahoor, M. H. (2016). AffectNet: A Database for Facial Expression, Valence, and Arousal Computing in the Wild. *IEEE Transactions on Affective Computing*, 10(1), 18–31.
12. Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.
13. Tang, Y. (2013). Deep learning using support vector machines. *arXiv preprint arXiv:1306.0239*.
14. Tzirakis, P., Trigeorgis, G., Nicolaou, M. A., Schuller, B. W., & Zafeiriou, S. (2017). End-to-End Multimodal Emotion Recognition using Deep Neural Networks. *IEEE Journal of Selected Topics in Signal Processing*, 11(8), 1301–1309.
15. Viola, P., & Jones, M. (2001). Rapid object detection using a boosted cascade of simple features. *CVPR*.
16. Zhang, K., Zhang, Z., Li, Z., & Qiao, Y. (2016). Joint Face Detection and Alignment using Multitask Cascaded Convolutional Networks. *IEEE Signal Processing Letters*, 23(10), 1499–1503.

17. Zhou, H., & Shi, J. (2017). Spatial attention CNN for emotion recognition. *ICCV Workshops*.
18. Zeng, Z., Pantic, M., Roisman, G. I., & Huang, T. S. (2009). A survey of affect recognition methods: Audio, visual, and spontaneous expressions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(1), 39–58.
