# Multi-label Movie Genre Detection from a Movie Poster Using Knowledge Transfer Learning

Kaushil Kundalia[1] · Yash Patel[1] · Manan Shah[2]

## Abstract

The task of predicting a movie genre from its poster can be very challenging owing to the high variability of movie posters. A novel approach for the generation of a multi-label movie genre prediction from its poster using neural networks that employ knowledge transfer learning has been proposed in this paper. This approach works on two fronts; one is aimed at creating a large, diverse and balanced dataset for movie genre prediction. The second front involves reframing the problem to simpler single-label multi-class classification and generating a multi-label multi-class prediction on a given movie poster as input. The experimental evaluation suggests that our approach generates a remarkable accuracy which is a result of a larger, evenly distributed dataset, simplifying the problem to a single-label multi-class classification problem and because of the use of knowledge transfer learning to extract higher-level feature.

**Keywords** Computer vision · Deep learning · Inception V3 · Transfer learning · Movie genre · Convolutional neural network

## Introduction

Computer vision has improved by leaps and bounds over the years, and with the advent of convolution neural networks, it has become much easier to achieve computer vision tasks such as object identification [13, 16]. In spite of all of these, there has not been much success in deciding the type of movies merely by looking at its poster with neural networks [7, 8].

The use of accurate movie genre classification has become a very important part of the recommended systems. Multimedia streaming services, like Amazon Prime Videos and others, show recommended movies to the users by showing posters of movie. However, the movie posters are highly varied and a poster might not always entirely justify its genre. (e.g. As shown in Fig. 1b, it becomes challenging to justify as an action movie by looking just at the poster without knowing another context.) This makes the user less interested. (e.g. if an action movie interested user might by deviated by looking at the poster of Fig. 1b.) Hence, it has become important for such recommendation system to also classify the movie based on its poster, so that only the attractive posters that truly justify the genre of a movie to the humans can be shown to its users. This classification is termed as automatic movie genre detection from its poster.

To address this issue, Netflix generates personalized artwork for the members that help them to find content to watch and enjoy maximizing member satisfaction and retention. It makes use of static images generated from the source video and creates raw artwork. Netflix then ranks these images that meet the aesthetics, creativity, and diversity in objects to represent content accurately. This creates a relevant and personalized artwork for each member based on their interest. Our proposed approach aims at addressing similar issues, wherein we categorize the images based on their aesthetics into various genres.

✉ Manan Shah
manan.shah@spt.pdpu.ac.in

Kaushil Kundalia
kaushil.kundalia@gmail.com

Yash Patel
yash9132@gmail.com

1 Department of Computer Engineering, Indus University, Ahmedabad, Gujarat, India

2 Department of Chemical Engineering, School of Technology, Pandit Deendayal Petroleum University, Gandhinagar, Gujarat, India

**(a)**                  **(b)**                  **(c)**                  **(d)**

Fig. 1 **a** Buster Keaton in The General, 1926 (category, Action). **b** From Russia with love, 1963 (category, Action). **c** Leon, 1994 (category, Action). **d** American Sniper, 2014 (category, Action). This

figure is a comparison about the high diversity the movie posters among a single category, Action). It can be seen that posters from the 1930s differ a lot from 2014. This adds to the intra-class variability

Since our approach uses only movie posters to categorize the genre of movies, it can be used to generate such personalized feedbacks for members based on their recommendations.

Computer vision tasks like image classification are becoming increasingly improvised owing to the advancements in convolutional neural network architecture since AlexNet and increasing computational power [9]. Convolutional neural networks (CNNs) have become increasingly popular in the field of computer vision due to their ability to generalize well as compared to a normal fully connected dense network [12]. CNNs follow deep feed-forward architecture. Apart from that, CNNs have a remarkably lesser number of parameters that allow them to generalize in a better manner. However, despite the advancements and the ability of CNNs, to generalize, automatic movie genre detection can be considered as one of the most challenging tasks of classification by deep neural network. This is because of the high level of intra-class variability in the visual representation of a poster, which makes it extremely difficult for a model to create any sort of fixed template of any particular class. This variability is furthermore intensified when considering a dataset that consists of movie posters over a large span of time. That is, movie posters change rapidly with changing time and trends. For example, posters in the 1920s were mostly hand painted which differs totally from current movie posters, e.g. as shown in Fig. 1, different movies belong to various cultural backgrounds, also contribute to this variability. All these scenarios make classification task even for humans quite challenging. In addition to this variability, multi-label classification makes it more challenging to achieve noteworthy accuracy.

Previous work by Chu and Guo [2] and Ivasic-Kos et al. [6] have addressed this multi-label classification by
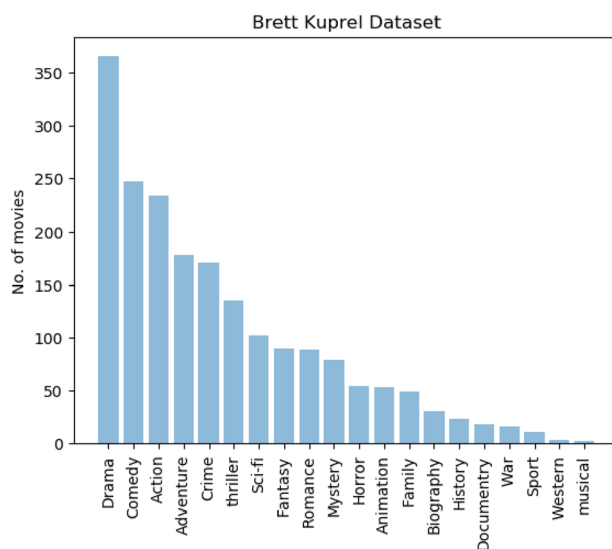


Fig. 2 Bar chart of Brett Kuprel's Dataset. This shows that the dataset is highly skewed and the model can potentially be biased towards Drama genre. This occurs because a movie can belong to multiple genres

training on a dataset with multi-labelled classes. This generates a highly skewed data (Figs. 2, 3). Further details about their approach have been explained in the following section. Krupel [10] extracted higher-level feature using CNNs; however, his dataset also suffered from the same problem of unbalanced data. To address this problem, we created our dataset, possibly the largest dataset on movie poster classification in our notice. Our approach consists of higher-level feature extraction using pertained InceptionV3 net, on the weight of ImageNet [9]. We stacked a series of fully connected dense layers on the InceptionV3 network. We trained our model on single-label class prediction. However, for prediction, the model predicts the top 3
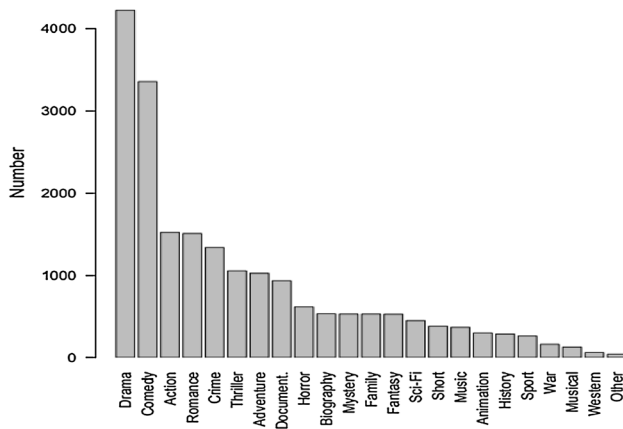
**Fig. 3** Distribution of dataset. This dataset is also highly skewed [2]

classes with the highest probability. Our approach has been further explained in detail in the upcoming sections.

## Related work

Huang et al. [4] in their work performed a single-label classification based on movie preview. Their approach classified movies among action, drama, and thriller category. Rather than making use of a single-frame movie poster, they extracted four features from a movie preview, namely average shot length, colour variance, motion content, and lighting key. Their paper shows strong evidence that combining visual cues with cinematic principles generates a good classifier.

Krupel [10], in his approach of classifying movie genres from the poster, made use of deep neural networks with auto-encoders for high-level feature extraction. His dataset consists of 5800 greyscale images of size $100 \times 100$ pixels from IMDb. He generated a prediction of a single-label class using a softmax activation function on the output layer. One of the major issues that Brett's model suffered was the unbalanced dataset. The summary is shown in Fig. 2.

Wehrmann and Barros [18] in his approach made use of movie trailer to predict multi-label genre by extracting higher-level features using CTT-MMC (Convolutional Through Time for Multi-Label Movie genre Classification). CTT-MC makes use of features generated from movie trailer frames. The CTT-MMC is an extension to the CTT model (Convolutional Through Time) that was originally created by Wehrmann and Barros [19] that is used to match the natural language description to the image context. The core concept behind CTT is to project the textual and visual features in the same embedding space. This makes it relate the naturally hierarchal concepts within image captions. CTT-MMC is in ultra-deep connection of

convolutional layers that have residual connections and makes use of a special convolutional layer to extract temporal information from image-based features before performing the mapping of movie trailers to genres. The first component of CTT-MMC is a convolutional network residual connections that are pre-trained on ImageNet and Places365 datasets [9, 21]. This component is responsible for extracting the higher-level features. CTT-MMC generates a temporal representation from the movie trailer. These temporal extracted features are further given as input to the CTT model that can generate temporal relationships among the temporal scenes. The output of the convolutional layers is then fed to a fully connected layer (FC) over a maxout activation function. The final predictions are computed by taking a logistic sigmoidal activation function of the generated output from the FC layers.

Rasheed and Shah [14] performed movie genre classification based on a movie preview. They implemented the classification on a decision tree-like fashion, wherein the initial classification is a binary classification between action and non-action movie. This classification is made based on visual disturbance. The action movies are further classified among Fire/Explosive and other classes. The non-action movies are further classified among Comedy, Horror and Drama/Other. The later classification is made based on audio, cinematic principles and colour. Authors work is one of the first steps towards automatic movie genre classification from a movie preview.

Zhou et al. [22] proposed an approach of scene categorization. They categorized movies into action, drama, horror, and comedy genre. They extracted a set of key frames from the movie trailer using shot boundary detection. They further used GIST, CENTRIST, and W-CENTRIST to extract higher-level features from these frames. This approach was a step further from using lower-level features. The extracted features were then clustered using K-means clustering to receive 100 feature clusters.

Chu and Guo [2] also created their own dataset of 8191 movie posters from IMDb with labels spanning over 23 classes of different movie genres. They proposed a novel approach of constructing a function 'F' that jointly considers a combination of visual representations generated by CNNs and the objects extracted from these posters using YOLO [15]. For the prediction of multi-labels, they proposed a threshold mechanism. A movie can belong to $i$th genre only if its probability surpasses theta threshold. To find the threshold, they used a grid search based on Matthews correlation coefficients to determine the best threshold in each dimension. To overcome the unbalanced class distribution of their dataset, they feature augmented images from the class having fewer instances. The summary of their dataset is shown in Fig. 3. Despite these

approaches, their model could generate a maximum accuracy of 18.73% [2].

Huang and Wang [5] employed a meta-heuristic optimization algorithm called self-adaptive harmony search (SAHS). They used a hybrid of visual and audio features extracted from a movie trailer and used a SVM with a one-vs-one setting for classification. They achieved a remarkable accuracy of 91.9%.

Ivasic-Kos et al. [6] initially generated a dataset of 6739 movies consisting of 18 genres (Action, Adventure, Animation, Comedy, Crime, Disaster, Documentary, Drama, Fantasy, History, Horror, Mystery, Romance, Science Fiction, Suspense, Thriller, War and Western). They picked 20 most popular movies from each year from 1990. Each movie could have one or more genres. The distribution of genre was highly skewed and some genre had a considerably fewer number of instances which was not enough to define an accurate classifier. To overcome this, they devised two approaches; firstly, they manually combined two or more similar genre into one, e.g. mystery and crime were combined with thriller. This reduced the number of classes to 11. Secondly, they discarded the remaining genres that did not fall among the previously

mentioned 11 classes. Their model is inspired by human which extracts lower-level features. They employed Naïve Bayes classifier, RAKEL, and ML-kNN [17, 20]. These models use a combination of features that have GIST features, dominant colour features (DC), local dominant colour features (DC1-DC5), and colour moments (CM).

## Our approach

One persistent conundrum that accounted for the low performance for these models was the existence of a highly unbalanced data as shown in Fig. 4. The statistical summary of various datasets is shown in Table 1. This is an inevitable situation while considering a dataset with multi-labels movie genre. This skewness exists because not all movies have an equal number of genres and also, not all classes are distributed evenly among all movies. For example, there are more number of movie releases that have genres of Action, Thriller, and Adventure combined than movies with Musical genre. This makes it impossible to generate an evenly distributed model that is not biased towards any particular class. To address this issue, we
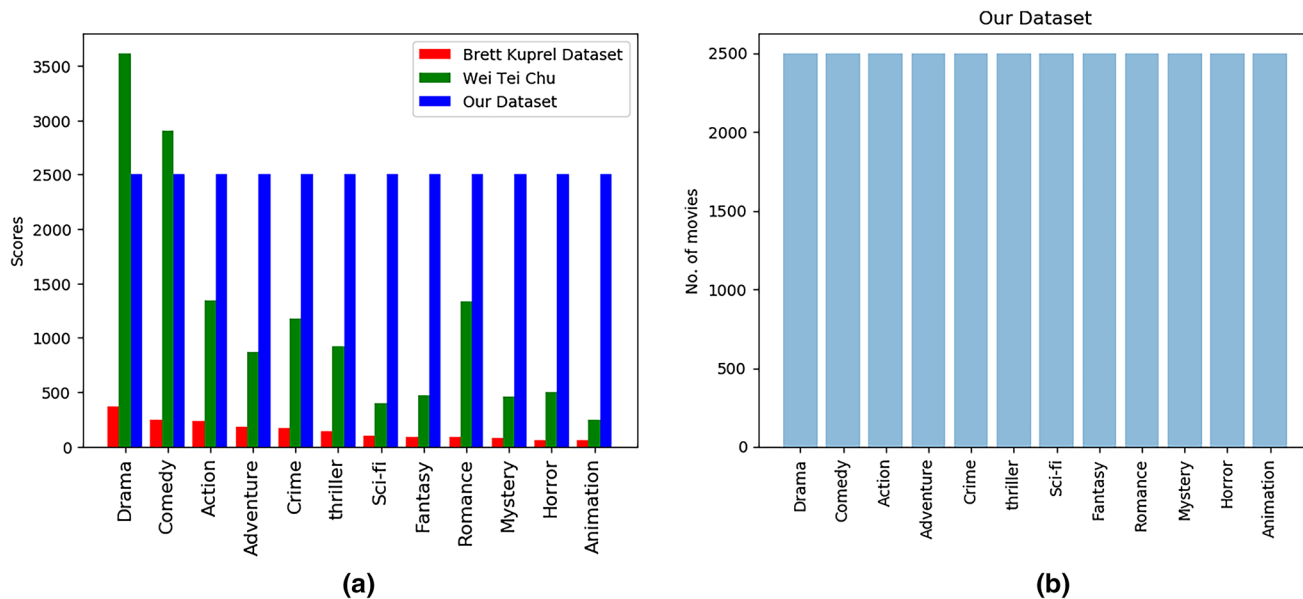


**Fig. 4** **a** Comparison by bar chart of occurrence of various genres among three datasets. As shows in the figure, our dataset is evenly distributed among all the 12 classes. **b** Bar chart of distribution of different classes among our dataset

**Table 1** Comparison of class distribution among various datasets prepared for movie genre classification from poster

| Research | Max per genre | Min per genre | Mean examples per genre | SD per genre |
|---|---|---|---|---|
| Marina Ivasic-Kos [5] | 417 | 1 | 4.6 | 14.87 |
| Wei-Ta Chu [2] | 3619 | 2 | 674.96 | 862.9 |
| Brett Kuprel [9] | 365 | 3 | 97.5 | 95 |
| Our dataset | 2500 | 2500 | 2500 | 0 |

Algorithm 1:Scraping for Comedy genre

```
Input: Desired number of movies N in multiples
of 50
Output: Posters classified into different folders
based on genre and a csv file which stores genre
g, movie title m , year y, and rating r.
e=N/50
c=0
U1=http://www.imdb.com/https://www.imdb.c
om/search/title/?genres=comedy&explore=title
_type,genres
while c<e do
      while c1<50 do
            p=getposter(U1)
            m=getTitle(U1)
            y=getYear(U1)
            g=getGenre(U1)
            r=getRating(U1)
            store p at folder g
            write m,y,g,r in csv file
      end while
      c1=0
      U1=getnextUrl(U1)
end while
```

**Fig. 5** Part of the algorithm used in scrapping for the dataset. Similar algorithm is modified for all other genres of movies

created a dataset that consists of only a single label associated with each movie. This way, we could eliminate the problem of unbalanced data by having each class of almost equal amount of occurrence. Although training a model on a dataset with single label would lead to a single-label prediction, to overcome this limitation, we used the following approach. Our model is defined on a single-labelled classification problem. However, to generate a multi-label prediction, given any movie poster, we define that each movie can belong to three genres (or classes). Hence, each predicted output generated by the model will be a set of three movie genres. The top 3 classes with maximum probability distribution will be returned as the output classes.

## Dataset

As mentioned above, to address the issue of unbalanced classes, we created a new dataset of 12 genres of movies with over 30,000 images. All the classes are evenly distributed, and each class has around 2500 images. This dataset is extracted from IMDb using a web crawler. The various classes include Action, Adventure, Animation, Comedy, Crime, Drama, Fantasy, Horror, Mystery, Romance, Sci-fic, and Thriller. The comparison of our dataset with various other existing dataset is shown in Fig. 4. Each image in our dataset is $960 \times 600$ pixels, with RGB channel.

This dataset has movie posters ranging from 1902 AD to 2019 AD. This covers a very large variability because of changing trends and techniques used in making movie posters.

To encourage advancements in the field of movie genre classification, we have made this dataset publicly available for everyone. This dataset probably might be the largest dataset made in our notice.
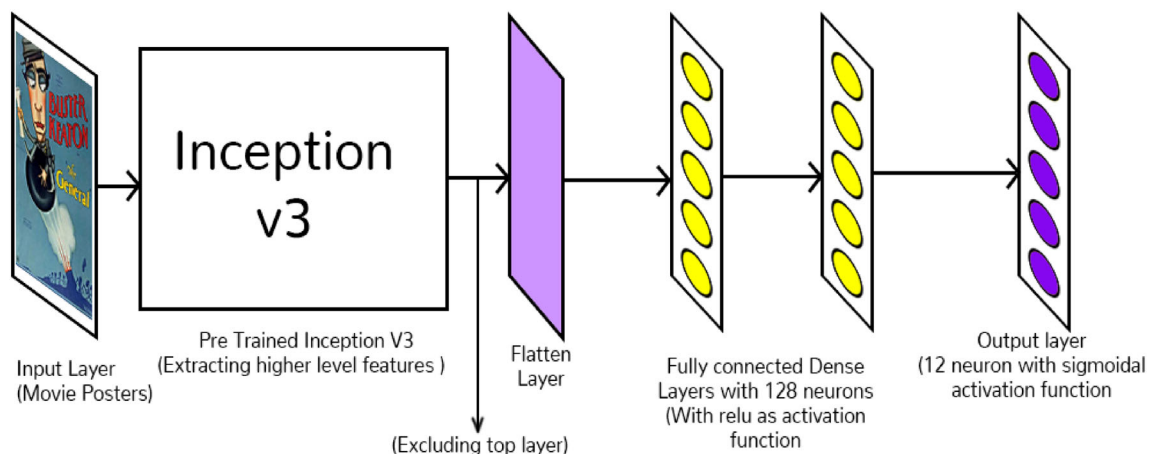


## Proposed Architecture

**Fig. 6** Proposed model architecture

(a)                    (b)                    (c)                    (d)

| Figure | a | b | c | d |
|---|---|---|---|---|
| Actual Movie Genre | Crime, Thriller | Comedy, Drama, Romance | Comedy, Romance | Action, Horror, Fantasy |
| Detected Movie Genre | ● Thriller 0.291 <br> ● Mystery 0.27159 <br> ● Crime 0.2706 | ● Comedy 0.3945 <br> ● Romance 00.2846 <br> ● Fantasy:0. 181 | ● Romance: 0.5475 <br> ● Comedy 0.410612 <br> ● Drama: 0.250 | ● Horror 0.3568 <br> ● Action 0.259 <br> ● Thriller 0.2465 |

Fig. 7 Movie poster examples which are correctly classified

## Experimental evaluation

### Transfer learning

One of the major contributors to the advancements in the field of deep learning is the ability to transfer knowledge. Transfer learning aims to provide a framework to utilize previously acquired knowledge to solve new problems much quickly and efficiently [11]. Transfer learning makes use of pre-trained models that are trained over a large dataset. These models are then repurposed according to the new problem definition. As the pre-trained models are trained over a huge dataset (e.g. inception is trained over ImageNet having about 14 million images with 1000 classes), it is safe to say that these models can generalize well. Because of its ability to generalize, transfer learning is highly used in computer vision. These models are trained for a large amount of time that requires high computational resources. Due to our limitation of computational power, we have employed the use of inception model to extract higher-level features in a movie poster (Fig. 5).

| Figure | a | b | c | d |
|---|---|---|---|---|
| Actual Movie Genre | Comedy, Drama | drama, thriller | Drama | Drama |
| Detected Movie Genre | ● Horror 0.311, ● Thriller 0.2474, ● Action 0.24602 | ● Horror 0.417277, ● Sci-fi 0.3136 ● Action 0.2195 | ● Thriller 0.3921, ● Mystery 0.2532, ● Horror 0.25088 | ● Horror 0.332, ● Action 0.2470 ● Thriller 0.2281 |
|  |  |  |  |  |

**Fig. 8** Movie poster examples which are classified incorrectly

## Training method and model architecture

To extract higher-level features from posters, authors used pre-trained model of InceptinV3 on the weights of ImageNet. The inception model consists of a series of convolutional layers that are used for feature extraction. The top (output) layer of the pre-trained model was removed, and we stacked a series of fully connected dense layers. These dense layers repurposed the model to generate a movie genre classifier. Our proposed model architecture is shown in Fig. 6. Each dense layer has 128 neurons and uses a rectifier linear (ReLu) activation function. The output layer is a dense layer with 12 output nodes. This layer has a sigmoid activation function.

Each instance in the dataset used for training the model has features shape of $200 \times 150 \times 3$ dimensions. Each

label instance has a dimension of $1 \times 12$, where 12 corresponds to the number of classes to be predicted.

The model was trained over a training set consisting of 30,000 images of equal class distribution. Each image is a 200 by 150 pixel, RGB image. We separated a validation set of 2450 images. We used stochastic gradient descent (SGD) optimizer and cross-entropy for loss function.

However, the limited processing resources made it difficult to train the model over the entire dataset at once. So we divided the training dataset into randomly shuffled six subsets with exclusive elements. The model was then trained individually over each such subset. We defined one epoch as the model training over the entire dataset once, instead of the model being trained over one subset. We used Keras API with Tensorflow backend using Python [1].

## Testing

Our dataset is a single-label prediction dataset. So we cannot evaluate the accuracy on a testing subset of that dataset. To overcome this challenge, we took 3450 random instances from Chu and Guo [2] dataset. They have designed a multi-label classification dataset, which is ideal for measuring the performance of our multi-label prediction problem. This created another inconsistency as their dataset is categories of a total of 25 genre (Action, Adventure, Animation, Biography, Comedy, Crime, Documentary, Drama, Family, Fantasy, History, Horror, Music, Musical, Mystery, N/A (other), News, Reality-TV, Romance, Sci-Fic, Short, Sports, Thriller, War and Western), and our dataset has only 12 genres (Action, Adventure, Animation, Comedy, Crime, Drama, Fantasy, Horror, Mystery, Romance, Sci-Fic, and Thriller). So we removed those genres from the testing set that were absent in our dataset. With this model, the achieved accuracy is 84.82% with a loss of 0.4892. Some of the movie genres predicted by the model are shown in Figs. 7, 8.

## Results

In Fig. 7, we have listed few randomly generated movies that were falsely predicted by the model; however, any human being would also have predicted the same genre without any context by merely looking that the same poster. This provides evidence that our model shows some resemblance to detecting movie posters as humans do. For example, movie posters with dominant dark-coloured pixels were mostly predicted as the horror genre. Figure 6 enlists some of the cases in which our algorithm predicted correctly.

In the future, we plan on improving the model by making use of the movie trailer, as inspired by Hui-Yu et al. [3] approach.

## Conclusion

From this research, it can be concluded that it is a challenge to predict a movie merely from its poster. This is because of a high level of variability and lack of pattern formation. However, this task can be made more efficient by having an equal distribution and relatively large dataset. Employing knowledge transfer learning to extract higher-level features can also be considerably useful. Further, it can also be argued that classifying a movie is a work of creativity and cannot be easily be quantified in numbers. Furthermore, movie posters are also designed in a way that reflects its current trend like from results; we can observe that movies releases recently have a difference in what poster conveys and what it actually is. This bias may be attributed to the fact that movies are made for commercial purposes, and hence posters are made in a way to attract larger crowds. Despite these flaws, this classification can give a general idea about the movie genre from its poster.

## Compliance with ethical standards

## References

1. Abadi M, Agarwal A, Barham P, Brevdo E, Chen Z, Citro C, Corrado GS, Davis A, Dean J, Devin M, Ghemawat S, Goodfellow IJ, Harp A, Irving G, Isard M, Jia Y, Józefowicz R, Kaiser L, Kudlur M, Levenberg J, Mane D, Monga R, Moore S, Murray DG, Olah C, Schuster M, Shlens J, Steiner B, Sutskever I, Talwar K, Tucker PA, Vanhoucke V, Vasudevan V, Viégas FB, Vinyals O, Warden P, Wattenberg, Wicke M, Yu Y, Zheng X (2016) Tensor flow: large-scale machine learning on heterogeneous distributed systems, pp 1–19
2. Chu W, Guo H (2017) Movie genre classification based on poster images with deep neural networks. In: MUSA2, proceedings of

the workshop on multimodal understanding of social, affective and subjective attributes. ACM, New York, NY, pp 39–45

3. Hui-Yu H, Weir-Sheng S, Wen-Hsing H (2008) A film classifier based on low-level visual features. J Multimed. https://doi.org/10.4304/jmm.3.3.26-33

4. Huang H, Shih W, Hsu WH (2008) Film classification based on low-level visual effect features. J Electron Imaging 17(2):023011

5. Huang YF, Wang SH (2012) Movie genre classification using SVM with audio and video features. In: Huang R, Ghorbani AA, Pasi G, Yamaguchi T, Yen NY, Jin B (eds) Active media technology. AMT 2012. Lecture Notes in Computer Science, 7669. Springer, Berlin, pp 1–10

6. Ivasic-Kos M, Pobar M, Mikec L (2014) Movie posters classification into genres based on low-level features. In: 37th international convention on information and communication technology, electronics and microelectronics (MIPRO), IEEE, pp 1198–1203

7. Jha K, Doshi A, Patel P, Shah M (2019) A comprehensive review on automation in agriculture using artificial intelligence. Artif Intell Agric 2:1–12

8. Kakkad V, Patel M, Shah M (2019) Biometric authentication and image encryption for image security in cloud framework. Multiscale Multidiscip Model Exp Des. https://doi.org/10.1007/s41939-019-00049-y

9. Krizhevsky A, Sutskever I, Hinton GE (2012) ImageNet classification with deep convolutional neural networks. In: Part of: Advances in neural information processing systems. NIPS 2012, vol 25, pp 1–9

10. Krupel B (2016) Judging a movie by its poster using deep learning. Stanford CS221, pp 1–5

11. Lu J, Behbood V, Hao P, Zuo H, Xue S, Zhang G (2015) Transfer learning using computational intelligence: a survey. Knowl-Based Syst 80:14–23

12. Nebauer C (1998) Evaluation of convolutional neural networks for visual recognition. IEEE Trans Neural Netw 9(4):685–696

13. Pandya R, Nadiadwala S, Shah R, Shah M (2019) Buildout of methodology for meticulous diagnosis of K-complex in EEG for aiding the detection of Alzheimer's by artificial intelligence.

Augment Human Res. https://doi.org/10.1007/s41133-019-0021-6

14. Rasheed Z, Shah M (2002) Movie genre classification by exploiting audio-visual features of previews. In: Object recognition supported by user interaction for service robots, vol 2, pp 1086–1089

15. Redmon J, Farhadi A (2017) YOLO9000: better, faster, stronger. In: IEEE conference on computer vision and pattern recognition (CVPR), Honolulu, HI, pp 6517–6525

16. Shah G, Shah A, Shah M (2019) Panacea of challenges in real-world application of big data analytics in healthcare sector. Data Inf Manag. https://doi.org/10.1007/s42488-019-00010-1

17. Tsoumakas G, Vlahavas I (2007) Random k-labelsets: an ensemble method for multilabel classification. In: Kok JN, Koronacki J, Mantaras RL, Matwin S, Mladenič D, Skowron A (eds) Machine learning: ECML 2007. ECML 2007. Lecture Notes in Computer Science, vol 4701. Springer, Berlin, pp 406–417

18. Wehrmann J, Barros RC (2017) Movie genre classification: a multi-label approach based on convolutions through time. Appl Soft Comput 61:973–982

19. Wehrmann J, Mattjie A, Barros RC (2018) Order embeddings and character-level convolutions for multimodal alignment. Pattern Recognit Lett 102:15–22

20. Zhang ML, Zhou ZH (2007) Ml-knn: a lazy learning approach to multi-label learning. Pattern Recognit 40(7):2038–2048

21. Zhou B, Lapedriza A, Xiao J, Torralba A, Oliva A (2014) Learning deep features for scene recognition using places database. In: Advances in neural information processing systems, pp 487–495

22. Zhou H, Hermans T, Karandikar AV, Rehg JM (2010) Movie genre classification via scene categorization. In: Proceedings of the international conference on multimedia, pp 747–750