100876

:

: 1)        ; 2)
（      、        ）        。

。

:        。

。

:        。

。

:        （FER）;        ;        ;

# Deep facial expression recognition: a survey

Li Shan    Deng Weihong

*School of Artificial Intelligence, Beijing University of Posts and Telecommunications, Beijing 100876, China*

**Abstract**: Facial expression is a powerful, natural, and universal signal for human beings to convey their emotional states and intentions. Numerous studies have been conducted on automatic facial expression analysis because of its practical importance in sociable robotics, medical treatment, driver fatigue surveillance, and many other human-computer interaction systems. Various facial expression recognition (FER) systems have been explored to encode expression information from facial representations in the field of computer vision and machine learning. Traditional methods typically use handcrafted features or shallow learning for FER. However, related studies have collected training samples from challenging real-world scenarios, which implicitly promote the transition of FER from laboratory-controlled to in-the-wild settings since 2013. Meanwhile, studies in various fields have increasingly used deep learning methods, which achieve state-of-the-art recognition accuracy and remarkably exceed the results of previous investigations due to considerably improved chip processing abilities (e. g., GPU units) and appropriately designed network architectures. Moreover, deep learning techniques are increasingly utilized to handle challenging factors for emotion recognition in the wild because of the effective training of facial expression data. The transition of facial expression recognition from being laboratory-controlled to challenging in-the-wild conditions and the recent success of deep learning techniques in various fields have promoted the use of deep neural

networks to learn discriminative representations for automatic FER. Recent deep FER systems generally focus on the following important issues. 1) Deep neural networks require a large amount of training data to avoid overfitting. However existing facial expression databases are insufficient for training common neural networks with deep architecture which achieve promising results in object recognition tasks. 2) Expression-unrelated variations are common in unconstrained facial expression scenarios such as illumination head pose and identity bias. These disturbances are nonlinearly confounded with facial expressions and therefore strengthen the requirement of deep networks to address the large intraclass variability and learn effective expression-specific representations. We provide a comprehensive review of deep FER including datasets and algorithms that provide insights into these intrinsic problems in this survey. First we introduce the background of fields of FER and summarize the development of available datasets widely used in the literature as well as FER algorithms in the past 10 years. Second we divide the FER system into two main categories according to feature representations namely static image and dynamic sequence FER. The feature representation in static-based methods is encoded with only spatial information from the current single image whereas dynamic-based methods consider temporal relations among contiguous frames in input facial expression sequences. On the basis of these two vision-based methods other modalities such as audio and physiological channels have also been used in multimodal sentiment analysis systems to assist in FER. Although pure expression recognition based on visible face images can achieve promising results incorporating it with other models into a high-level framework can provide complementary information and further enhance the robustness. We introduce existing novel deep neural networks and related training strategies which are designed for FER based on both static and dynamic image sequences and discuss their advantages and limitations in state-of-the-art deep FER. Competitive performance and experimental comparisons of these deep FER systems in widely used benchmarks are also summarized. We then discuss relative advantages and disadvantages of these different types of methods with respect to two open issues (data size requirement and expression-unrelated variations) and other focuses (computation efficiency performance and network training difficulty) . Finally we review and summarize the following challenges in this field and future directions for the design of robust deep FER systems. 1) Lacking training data in terms of both quantity and quality is a main challenge in deep FER systems. Abundant sample images with diverse head poses and occlusions as well as precise face attribute labels including expression age gender and ethnicity are crucial for practical applications. The crowdsourcing model under the guidance of expert annotators is a reasonable approach for massive annotations. 2) Data bias and inconsistent annotations are very common among different facial expression datasets due to various collecting conditions and the subjectiveness of annotating. Furthermore the FER performance fails to improve when training data is enlarged by directly merging multiple datasets due to inconsistent expression annotations. Cross-database performance is an important evaluation criterion of generalizability and practicability of FER systems. Deep domain adaption and knowledge distillation are promising trends to address this bias. 3) Another common issue is imbalanced class distribution in facial expression due to the practicality of sample acquirement. One solution is to resample and balance the class distribution on the basis of the number of samples for each class during the preprocessing stage using data augmentation and synthesis. Another alternative is to develop a cost-sensitive loss layer for reweighting during network work training. 4) Although FER within the categorical model has been extensively investigated the definition of prototypical expressions covers only a small portion of specific categories and cannot capture the full repertoire of expressive behavior for realistic interactions. Incorporating other affective models such as FACS (facial action coding system) and dimensional models can facilitate the recognition of facial expressions and allow them to learn expression-discriminative representations. 5) Human expressive behavior in realistic applications involves encoding from different perspectives with facial expressions as only one modality. Although pure expression recognition based on visible face images can achieve promising results incorporating it with other models into a high-level framework can provide complementary information and further enhance the robustness. For example the fusion of other modalities such as the audio information infrared images and depth information from 3D face models and physiological data has become a promising research direction due to the large complementarity of facial expressions and the good application value of human-computer interaction (HCI) applications.

**Key words**: facial expression recognition (FER) ; real world; deep learning; survey

**0**

（Darwin　Prodger　1998；Tian　2001）。

；

Ekman　6

Ekman　Friesen（1971）　　　　　　　　（Ekman　1994）　6　：　、　、　、　、　。　　　　　6

：

（local binary pattern LBP）（Shan　2009）、（local binary pattern from three orthogonal planes LBP-TOP）（Zhao　Pietikainen 2007）、（nonnegative matrix factorization NMF）（Zhi　2011）（Zhong　2012）。2013　FER2013（the Facial Expression Recognition 2013）（Goodfellow　2013）EmotiW（Dhall　2015 2016 2017）

、

（　GPU　）

（Krizhevsky　2012；Simonyan　Zisserman 2014b；Szegedy　2015）。

1　　　　2007

DAE（deep autoencoder）、LP（locality-preserving）loss　IACNN（identity-aware CNN）。



图1

Fig. 1　The evolution of facial expression recognition and facial expression dataset

（Valstar　2012）

、

。Levi  Hassner( 2015)

LBP

Zhang  ( 2016)

( scale-invariant feature transform  SIFT)  ( Lowe

1999)

。

( Zeng  2018; Luo  2017)  、

、

。

。

( Chen  2018a)  3  ( region of inter-

est  ROI)  、

3  。

( Mavani  2017; Wu  Lin 2018)

。

1. 1. 2

( Ciregan  2012)。

: 1)

; 2)

。

1

。  ( Kim  2016)

。  、

( Kim

2015; Pons  Masip 2018b)。

2

:  。

( Bargal  2016; Liu  2016)。

3  :  、

。

。 Kahou  ( 2013)

。 Yu  Zhang( 2015)

。 Kim  ( 2015)

。 Pons  Masip

( 2018b)

。

1. 1. 3

。

。

、

。

Reed  ( 2014)  ( dis-

entangling Boltzmann machines  disBM)

。

( Devries  2014) 、  ( Pons

Masip 2018a)  ( Zhang  2017)

。

1. 1. 4

。

。

。 Lyu  ( 2014)

( deep belief network  DBN)

( stacked autoencoder  SAE)  。 Rifai

( 2012)  ( con-

tractive convolutional network  CCNET)

( con-

tractive autoencoder  CAE)

。

Liu  ( 2014)

( boosted deep belief network  BD-

BN)

、  。

。

。

1. 1. 5 　　　　　　　　　　　　　　　　　　　　　　　　1
　　　　　　　　　　　　　　　　　　　　（conditional GAN　cGAN）
　　　　　　　（generative adversarial networks　　　　　2
GAN）
　、　　　　　　　　　　　　　　　　。Yang　　（2018a）
　　　　　　　　　　　　GAN　　　（de-expression residue learning DeRL）

　　　　　　。　　　　　　　　　　　　　　　　　　　　。
　　　　　　　　　　　　　　Lai　Lai
（2018）　　　　　GAN　　　　　　　　　　　　　　　　　。
　　　　　　　　　　　　　　　　　1. 1. 6
　　　　　　　　　　　　　　　1
　　　　。Zhang　　（2018a）　　　　　　　　　　　　　　　。
　　GAN　　　　　　　　　　　　　　　　　　　　　　　　　。
　　　　　。
Yang　　（2018b）

**1**

**Table 1　Performance summary of representative methods for static-based deep
facial expression recognition on the most widely evaluated datasets**

| | | /% |
|---|---|---|
| CK + （Lucey　2010） | Ding　（2017） | 6　：（98. 6）8　：（96. 8） |
| | Zeng　（2018） | 7　：95. 79（93. 78）8　：89. 84（86. 82） |
| | Meng　（2017） | 7　：95. 37 （95. 51） |
| | Liu　（2017） | 7　：97. 1 （96. 1） |
| | Yang　（2018a） | 7　：97. 30 （96. 57） |
| | Zhang　（2018b） | 6　：98. 9 |
| MMI（Pantic　2005） | Liu　（2017） | 6　：78. 53 （73. 50） |
| | Li　（2017） | 6　：78. 46 |
| | Yang　（2018a） | 6　：73. 23 （72. 67） |
| | Liu　（2019） | 6　：81. 13 （79. 33） |
| FER2013（Goodfellow　2013） | Guo　（2016） | ：71. 33 |
| | Kim　（2016） | ：73. 73 |
| | Georgescu　（2019） | ：75. 42 |
| SFEW2. 0（Dhall　2015） | Li　（2017） | ：51. 05 |
| | Ding　（2017） | ：55. 15 （46. 6） |
| | Liu　（2017） | ：54. 19 （47. 97） |
| | Meng　（2017） | ：50. 98 （42. 57） |

：1)　　　　　　　　　　　　　　　　　　　　。2)7 ：　、　、　、　、　　　　　；8 ：　、
　、　、　、　、　。CK +　the extended Cohn-Kanade database　MMI　Maja Pantic　Michel Valstar Initiative。

。　　　　　　　　　。

。

1

（Jung　　2015）
CNN　　　（Yan　　2016）。

（Kim　　2017）。
（Zhang　　2017）

。

。　　　　　　　　　　。

。　　　　　　　　　　1. 2. 4

。　　　　　　　　　　　　　　　CNN
LSTM　　　　　　　　／　　　　　Donahue
（2015）　　　　　　　　　　　　　　CNN
。　　　　　　　　　　（　　　　　） LSTM

**1. 2**

。

。

。

1. 2. 1

（recurrent neural network
RNN）　　　　　　　　　。
（long short-term memory　LSTM）

1. 2. 2　3
RNN　　　　　　　　（convolutional neu-
ral networks　CNN）
3　　　　（3D convolutional neural network　C3D）

CNN　　2　　　C3D
3

1. 2. 3

（　　、　　　）

（Kim　　2019; Fan　　2016; Vielzeuf
2017; Jain　　2017）。　　　LSTM　CNN
Kankanamge　　（2017）　　　CNN
LSTM

。Ouyang
（2017）　　　　　　　　　　ResNet-LSTM
CNN　　　　　　　　　LSTMs
。CNN　Baccouche
（2012）

;　　　　　LSTM
LSTM Hasani　　Mahoor（2017）
（conditional random fields　CRFs）

。

1. 2. 5

Simonyan　　Zisserman
（2014a）

。

。Sun　　（2019）

。Zhang （2017） PHRNN 。RNN
（part-based hierarchical bidirectional recurrent neural LSTM
network） MSCNN（multisignal CNN） 。

／ 、 ／ ／ 。 3

。 3

。 Jung

（2015） 。

。

。

1.2.6 。

2 。

。 。

**2**

**Table 2 Performances of representative methods for dynamic-based deep
facial expression recognition on the most widely evaluated datasets**

| | | | /% |
|---|---|---|---|
| | Sun （2019） | | 6 ：97.28 |
| | Kumawat （2019） | 3 | 7 ：97.38（96.65） |
| CK＋（Lucey 2010） | Jung （2015） | | 7 ：97.25（95.22） |
| | Zhang （2017） | | 7 ：98.50（97.78） |
| | Hasani Mahoor（2017） | | 6 ：77.50（74.50） |
| MMI（Pantic 2005） | Zhang （2017） | | 6 ：81.18（79.30） |
| | Wang （2020） | | 6 ：82.21 |
| | Sun （2019） | | 6 ：91.46 |
| | Jung （2015） | | 6 ：81.46（81.49） |
| Oulu-CAISA（Zhao 2011） | Kumawat （2019） | 3 | 6 ：82.41（82.41） |
| | Zhang （2017） | | 6 ：86.25（86.25） |
| | Yan （2016） | VGG16-LSTM | 7 ：44.46 |
| | Yan （2016） | | 7 ：37.37 |
| AFEW 6.0（Dhall 2016） | Fan （2016） | VGG16-LSTM | ：45.43（38.96） |
| | Fan （2016） | 3 | ：39.69（38.55） |
| | Yan （2016） | | ：56.66（40.81） |
| | Fan （2016） | | ：59.02（44.94） |
| | Ouyang （2017） | VGG-LSTM | ：47.4 |
| AFEW 7.0（Dhall 2017） | Ouyang （2017） | 3 | ：35.2 |
| | Vielzeuf （2017） | VGG16-LSTM | ：48.6 |
| | Vielzeuf （2017） | | ：58.81（43.23） |

： 。Oulu-CAISA Oulu-CAISA facial expression database。

中国图象图形学报
JOURNAL OF IMAGE AND GRAPHICS

。

。

## 2

**2. 1**

。

。　　　　　　　　、

、　　　　。

。

。

。

。Benitez-Quiroz　（2016）

。2017

flicker　　　3

40

。

RAF-DB（Real-world affective face data-base）（Li　2017; Li　Deng　2018）。

EmotiNet（Benitez-Quiroz　2016）　AffectNet（Mol-lahosseini　2019）

。

。

**2. 2**

7

:

。　　　　　　　、

。

:　　　　　　　　　　（facial action coding system FACS）（Ekman　Rosenberg　1997）

;

（Yan　2013　2014）。
（Gunes　Schuller　2013; Russell　1980）

。　　　　　　　　　　　　Du

（2014）

。

。RAF-DB（Li　Deng　2019b）
RAF-ML（Real-world affective face multiLabel）（Li Deng　2019a）　　　　　7　　　　　　、
12　　　　　　　30　　　　　。Lu
（2018）　　　　　　　　　　　　　Word-net

。

:

;

。

**2. 3**

。

。

（Li　Deng　2020a）

。

（Wei
2018; Li　Deng　2020a）。

。

:

、

。

（

）　　　　　　　　。

。

**2. 4**

、3

—

。

。

、3

。

（remote photoplethysmography rPPG）
RhythmNet（Niu　　2019）

。

（Li　　2019c）

、

。

（Li　　2019b）

。

（Li　　2019a）

。

（Wang
2020; Li　　2019b）。3　　　　　　　（Chen
2018b）

（Zheng　　2018）、　　　　（Wang　　2018）
（Liu　　2020）
。Li　　Deng（2020b）

。

**3**

2007

（**References**）

Baccouche M　Mamalet F　Wolf C　Garcia C and Baskurt A. 2012. Spa‐
tio‐temporal convolutional sparse auto‐encoder for sequence classifi‐
cation//Proceedings of the British Machine Vision Conference. Sur‐
rey: BMVA Press: 1‐12　DOI: 10. 5244/C. 26. 124

Bargal S A　Barsoum E　Ferrer C C and Zhang C. 2016. Emotion recog‐
nition in the wild from videos using images//Proceedings of the 18th
ACM International Conference on Multimodal Interaction. New
York: ACM: 433‐436　DOI: 10. 1145/2993148. 2997627

Benitez‐Quiroz C F　Srinivasan R and Martinez A M. 2016. EmotioNet:
an accurate　real‐time algorithm for the automatic annotation of a
million facial expressions in the wild//Proceedings of 2016 IEEE
Conference on Computer Vision and Pattern Recognition. Las Ve‐
gas: IEEE: 5562‐5570　DOI: 10. 1109/CVPR. 2016. 600

Chen L F　Zhou M T　Su W J　Wu M　She J H and Hirota K. 2018a.
Softmax regression based deep sparse autoencoder network for facial
emotion recognition in human‐robot interaction. Information Sci‐
ences 428: 49‐61　DOI: 10. 1016/j. ins. 2017. 10. 044

Chen Z X　Huang D　Wang Y H and Chen L M. 2018b. Fast and light
manifold CNN based 3D facial expression recognition across pose
variations//Proceedings of the 26th ACM International Conference
on Multimedia. New York: ACM: 229‐238　DOI: 10. 1145/
3240508. 3240568

Ciregan D　Meier U and Schmidhuber J. 2012. Multi‐column deep neu‐
ral networks for image classification//Proceedings of 2012 IEEE
Conference on Computer Vision and Pattern Recognition. Provi‐

dence: IEEE: 3642-3649  DOI: 10. 1109/CVPR. 2012. 6248110

Darwin C and Prodger P. 1998. The Expression of the Emotions in Man and Animals. Oxford: Oxford University Press  1998

Devries T  Biswaranjan K and Taylor G W. 2014. Multi-task learning of facial landmarks and expression//Proceedings of 2014 Canadian Conference on Computer and Robot Vision. Montreal: IEEE: 98-103  DOI: 10. 1109/CRV. 2014. 21

Dhall A  Goecke R  Ghosh S  Joshi J  Hoey J and Gedeon T. 2017. From individual to group-level emotion recognition: EmotiW 5. 0//Proceedings of the 19th ACM International Conference on Multimodal Interaction. New York: ACM: 524-528  DOI: 10. 1145/3136755. 3143004

Dhall A  Goecke R  Joshi J  Hoey J and Gedeon T. 2016. EmotiW 2016: video and group-level emotion recognition challenges//Proceedings of the 18th ACM International Conference on Multimodal Interaction. New York: ACM: 427-432  DOI: 10. 1145/2993148. 2997638

Dhall A  Murthy O V R  Goecke R  Joshi J and Gedeon T. 2015. Video and image based emotion recognition challenges in the wild: emotiW 2015//Proceedings of 2015 ACM on International Conference on Multimodal Interaction. New York: ACM: 423-426  DOI: 10. 1145/2818346. 2829994

Ding H  Zhou S K and Chellappa R. 2017. facenet2expnet: regularizing a deep face recognition net for expression recognition//Proceedings of the 12th IEEE International Conference on Automatic Face and Gesture Recognition. Washington: IEEE: 118-126  DOI: 10. 1109/FG. 2017. 23

Donahue J  Hendricks L A  Guadarrama S  Rohrbach M  Venugopalan S  Darrell T and Saenko K. 2015. Long-term recurrent convolutional networks for visual recognition and description//Proceedings of 2015 IEEE Conference on Computer Vision and Pattern Recognition. Boston: IEEE: 2625-2634  DOI: 10. 1109/CVPR. 2015. 7298878

Du S C  Tao Y and Martinez A M. 2014. Compound facial expressions of emotion. Proceedings of the National Academy of Sciences of the United States of America  111(15): E1454-E1462  DOI: 10. 1073/pnas. 1322355111

Ekman P. 1994. Strong evidence for universals in facial expressions: a reply to Russell's mistaken critique. Psychological Bulletin  115(2): 268-287  DOI: 10. 1037/0033-2909. 115. 2. 268

Ekman P and Friesen W V. 1971. Constants across cultures in the face and emotion. Journal of Personality and Social Psychology  17(2): 124-129  DOI: 10. 1037/h0030377

Ekman P and Rosenberg E L. 1997. What the Face Reveals: Basic and Applied Studies of Spontaneous Expression using the Facial Action Coding System (FACS). New York: Oxford University Press  1997

Fan Y  Lu X J  Li D and Liu Y L. 2016. Video-based emotion recognition using CNN-RNN and C3D hybrid networks//Proceedings of the 18th ACM International Conference on Multimodal Interaction. New

York: ACM: 445-450  DOI: 10. 1145/2993148. 2997632

Georgescu M I  Ionescu R T and Popescu M. 2019. Local learning with deep and handcrafted features for facial expression recognition. IEEE Access  7: 64827-64836  DOI: 10. 1109/ACCESS. 2019. 2917266

Goodfellow I J  Erhan D  Carrier P L  Courville A  Mirza M  Hamner B  Cukierski W  Tang Y C  Thaler D  Lee D H  Zhou Y B  Ramaiah C  Feng FX  Li R F  Wang X J  Athanasakis D  Shawe-Taylor J  Milakov M  Park J  Ionescu R  Popescu M  Grozea C  Bergstra J  Xie J J  Romaszko L  Xu B  Chuang Z and Bengio Yoshua. 2013. Challenges in representation learning: a report on three machine learning contests//Proceedings of the 20th International Conference on Neural Information Processing. Daegu: Springer: 117-124  DOI: 10. 1007/978-3-642-42051-1_16

Gunes H and Schuller B. 2013. Categorical and dimensional affect analysis in continuous input: current trends and future directions. Image and Vision Computing  31(2): 120-136  DOI: 10. 1016/j. imavis. 2012. 06. 016

Guo Y N  Tao D P  Yu J  Xiong H  Li Y T and Tao D C. 2016. Deep neural networks with relativity learning for facial expression recognition//Proceedings of 2016 IEEE International Conference on Multimedia and Expo Workshops (ICMEW). Seattle: IEEE: 1-6  DOI: 10. 1109/ICMEW. 2016. 7574736

Hasani B and Mahoor M H. 2017. Spatio-temporal facial expression recognition using convolutional neural networks and conditional random fields//Proceedings of the IEEE International Conference on Automatic Face and Gesture Recognition. Washington: IEEE: 790-795  DOI: 10. 1109/FG. 2017. 99

Jain D K  Zhang Z and Huang K Q. 2017. Multi angle optimal pattern-based deep learning for automatic facial expression recognition. EB/OL. 2020-05-22. https://www. sciencedirect. com/science/article/pii/S0167865517302313

Jung H  Lee S  Yim J  Park S and Kim J. 2015. Joint fine-tuning in deep neural networks for facial expression recognition//Proceedings of 2015 IEEE International Conference on Computer Vision. Santiago: IEEE: 2983-2991  DOI: 10. 1109/ICCV. 2015. 341

Kahou S E  Pal C  Bouthillier X  Froumenty P  Gülçehre Ç  Memisevic V  Vincent P  Courville A  Bengio Y  Ferrari R C  Mirza M  Jean S  Carrier P L  Dauphin Y  Boulanger-Lewandowski N  Aggarwal A  Zumer J  Lamblin P  Raymond J P  Desjardins G  Pascanu R  Warde-Farley D  Torabi A  Sharma A  Bengio E  Côte M  Konda K R and Wu Z Z. 2013. Combining modality specific deep neural networks for emotion recognition in video//Proceedings of the 15th ACM on International Conference on Multimodal Interaction. New York: ACM: 543-550  DOI: 10. 1145/2522848. 2531745

Kankanamge S  Fookes C and Sridharan S. 2017. Facial analysis in the wild with LSTM networks//Proceedings of 2017 IEEE International Conference on Image Processing. Beijing: IEEE: 1052-1056  DOI: 10. 1109/ICIP. 2017. 8296442

Kaya H　Gürpınar F and Salah A A. 2017. Video-based emotion recognition in the wild using deep transfer learning and score fusion. Image and Vision Computing　65: 66-75　DOI: 10.1016/j. imavis. 2017. 01. 012

Kim B K　Dong S Y　Roh J　Kim G and Lee S Y. 2016. Fusing aligned and non-aligned face information for automatic affect recognition in the wild: a deep learning approach//Proceedings of 2016 IEEE Conference on Computer Vision and Pattern Recognition Workshops. Las Vegas: IEEE: 1499-1508　DOI: 10. 1109/CVPRW. 2016. 187

Kim B K　Lee H　Roh J and Lee S Y. 2015. Hierarchical committee of deep CNNs with exponentially-weighted decision fusion for static facial expression recognition//Proceedings of 2015 ACM on International Conference on Multimodal Interaction. New York: ACM: 427-434　DOI: 10. 1145/2818346. 2830590

Kim D H　Baddar W J　Jang J and Ro Y M. 2019. Multi-objective based spatio-temporal feature representation learning robust to expression intensity variations for facial expression recognition. IEEE Transactions on Affective Computing　10(2): 223-236　DOI: 10. 1109/TAFFC. 2017. 2695999

Kim D H　Lee M K　Choi D Y and Song B C. 2017. Multi-modal emotion recognition using semi-supervised learning and multiple neural networks in the wild//Proceedings of the 19th ACM International Conference on Multimodal Interaction. New York: ACM: 529-535　DOI: 10. 1145/3136755. 3143005

Knyazev B　Shvetsov R　Efremova N and Kuharenko A. 2017. Convolutional neural networks pretrained on large face recognition datasets for emotion classification from video. EB/OL. 2020-05-31　https: //arxiv. org/pdf/1711. 04598. pdf

Krizhevsky A　Sutskever I and Hinton G E. 2012. ImageNet classification with deep convolutional neural networks//Proceedings of the 25th International Conference on Neural Information Processing Systems. Red Hook: ACM: 1097-1105

Kumawat S　Verma M and Raman S. 2019. LBVCNN: local binary volume convolutional neural network for facial expression recognition from image sequences//Proceedings of 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW). Long Beach: IEEE: 207-216　DOI: 10. 1109/CVPRW. 2019. 00030

Lai Y H and Lai S H. 2018. Emotion-preserving representation learning via generative adversarial network for multi-view facial expression recognition//Proceedings of the 13th IEEE International Conference on Automatic Face and Gesture Recognition. Xi'an: IEEE: 263-270　DOI: 10. 1109/FG. 2018. 00046

Levi G and Hassner T. 2015. Emotion recognition in the wild via convolutional neural networks and mapped binary patterns//Proceedings of 2015 ACM on International Conference on Multimodal Interaction. New York: ACM: 503-510　DOI: 10. 1145/2818346. 2830587

Li S and Deng W H. 2018. Deep emotion transfer network for cross-database facial expression recognition//Proceedings of the 24th International Conference on Pattern Recognition. Beijing: IEEE: 3092-3099　DOI: 10. 1109/ICPR. 2018. 8545284

Li S and Deng W H. 2019a. Blended emotion in-the-wild: multi-label facial expression recognition using crowdsourced annotations and deep locality feature learning. International Journal of Computer Vision　127(6): 884-906　DOI: 10. 1007/s11263-018-1131-1

Li S and Deng W H. 2019b. Reliable crowdsourcing and deep locality-preserving learning for unconstrained facial expression recognition. IEEE Transactions on Image Processing　28(1): 356-370　DOI: 10. 1109/TIP. 2018. 2868382

Li S and Deng W H. 2020a. A deeper look at facial expression dataset bias. EB/OL. 2020-05-01. https: //arxiv. org/pdf/1904. 11150. pdf

Li S and Deng W H. 2020b. Deep facial expression recognition: a survey. EB/OL. 2020-05-31. https: //ieeexplore. ieee. org/document/9039580

Li S　Deng W H and Du J P. 2017. Reliable crowdsourcing and deep locality-preserving learning for expression recognition in the wild//Proceedings of 2017 IEEE Conference on Computer Vision and Pattern Recognition. Honolulu: IEEE: 2584-2593　DOI: 10. 1109/CVPR. 2017. 277

Li Y　Zeng J B　Shan S G and Chen X L. 2019a. Occlusion aware facial expression recognition using CNN with attention mechanism. IEEE Transactions on Image Processing　28(5): 2439-2450　DOI: 10. 1109/TIP. 2018. 2886767

Li Y　Zeng J B　Shan S G and Chen X L. 2019b. Self-supervised representation learning from videos for facial action unit detection//Proceedings of 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition. Long Beach: IEEE: 10924-10933　DOI: 10. 1109/CVPR. 2019. 01118

Li Y　Zheng W M　Wang L　Zong Y and Cuio Z. 2019c. From regional to global brain: a novel hierarchical spatial-temporal neural network model for EEG emotion recognition. IEEE Transactions on Affective Computing: #2922912　DOI: 10. 1109/TAFFC. 2019. 2922912

Liu K　Zhang M M and Pan Z G. 2016. Facial expression recognition with CNN ensemble//Proceedings of 2016 International Conference on Cyberworlds. Chongqing: IEEE: 163-166　DOI: 10. 1109/CW. 2016. 34

Liu P　Han S Z　Meng Z B　and Tong Y. 2014. Facial expression recognition via a boosted deep belief network//Proceedings of 2014 IEEE Conference on Computer Vision and Pattern Recognition. Columbus: IEEE: 1805-1812　DOI: 10. 1109/CVPR. 2014. 233

Liu P　Wei Y C　Meng Z B　Deng W H　Zhou J T and Yang Y. 2020. Omni-supervised facial expression recognition: a simple baseline. EB/OL. 2020-05-22. https: //arxiv. org/pdf/2005. 08551. pdf

Liu X F　Kumar B V K V　Jia P and You J. 2019. Hard negative generation for identity-disentangled facial expression recognition. Pattern

Recognition　88: 1-12　DOI: 10. 1016/j. patcog. 2018. 11. 001

Liu X F　Kumar B V K V　You J and Jia P. 2017. Adaptive deep metric learning for identity-aware facial expression recognition//Proceedings of 2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops. Honolulu: IEEE: 522-531　DOI: 10. 1109/CVPRW. 2017. 79

Lowe D G. 1999. Object recognition from local scale-invariant features//Proceedings of the International Conference on Computer Vision. Washington: ACM: 1150-1157

Lu Z J　Zeng J B　Shan S G and Chen X L. 2018. Zero-shot facial expression recognition with multi-label label propagation//Proceedings of the 14th Asian Conference on Computer Vision. Perth: Springer: 19-34　DOI: 10. 1007/978-3-030-20893-6_2

Lucey P　Cohn J F　Kanade T　Saragih J　Ambadar Z and Matthews I. 2010. The extended Cohn-Kanade dataset (CK +): a complete dataset for action unit and emotion-specified expression//Proceedings of 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition-Workshops. San Francisco: IEEE: 94-101　DOI: 10. 1109/CVPRW. 2010. 5543262

Luo Z J　Chen J H　Takiguchi T and Ariki Y. 2017. Facial expression recognition with deep age//Proceedings of 2017 IEEE International Conference on Multimedia and Expo Workshops (ICMEW). Hong Kong　China: IEEE: 657-662　DOI: 10. 1109/ICMEW. 2017. 8026251

Lv Y D　Feng Z Y and Xu C. 2014. Facial expression recognition via deep learning//Proceedings of 2014 International Conference on Smart Computing. Hong Kong　China: IEEE: 303-308　DOI: 10. 1109/SMARTCOMP. 2014. 7043872

Mavani V　Raman S and Miyapuram K P. 2017. Facial expression recognition using visual saliency and deep learning//Proceedings of 2017 IEEE International Conference on Computer Vision. Venice: IEEE: 2783-2788　DOI: 10. 1109/ICCVW. 2017. 327

Meng Z B　Liu P　Cai J　Han S Z and Tong Y. 2017. Identity-Aware Convolutional Neural Network for Facial Expression Recognition//Proceedings of the 12th IEEE International Conference on Automatic Face and Gesture Recognition. Washington: IEEE: 558-565　DOI: 10. 1109/FG. 2017. 140

Mollahosseini A　Hasani B and Mahoor M H. 2019. AffectNet: a database for facial expression　valence　and arousal computing in the wild. IEEE Transactions on Affective Computing　10(1): 18-31　DOI: 10. 1109/TAFFC. 2017. 2740923

Ng H W　Nguyen V D　Vonikakis V and Winkler S. 2015. Deep learning for emotion recognition on small datasets using transfer learning//Proceedings of the 17th ACM International Conference on Multimodal Interaction. New York: ACM: 443-449　DOI: 10. 1145/2818346. 2830593

Ng H W and Winkler S. 2014. A data-driven approach to cleaning large face datasets//Proceedings of 2014 IEEE International Conference on Image Processing. Paris: IEEE: 343-347　DOI: 10. 1109/ICIP. 2014. 7025068

Niu X S　Shan S G　Han H and Chen X L. 2019. RhythmNet: end-to-end heart rate estimation from face via spatial-temporal representation. IEEE Transactions on Image Processing　29: 2409-2423　DOI: 10. 1109/TIP. 2019. 2947204

Ouyang X　Kawaai S　Goh E G　Shen S M　Ding W　Ming H P and Huang D Y. 2017. Audio-visual emotion recognition using deep transfer learning and multiple temporal models//Proceedings of the 19th ACM International Conference on Multimodal Interaction. New York: ACM: 577-582　DOI: 10. 1145/3136755. 3143012

Pantic M　Valstar M　Rademaker R and Maat L. 2005. Web-based database for facial expression analysis//Proceedings of 2015 IEEE International Conference on Multimedia and Expo. Amsterdam: IEEE: #1521424　DOI: 10. 1109/ICME. 2005. 1521424

Parkhi O M　Vedaldi A and Zisserman A. 2015. Deep face recognition//Proceedings of the British Machine Vision Conference.　s. l. : BMVA Press: 41. 1-41. 12　DOI: 10. 5244/c. 29. 41

Pons G and Masip D. 2018a. Multi-task　multi-label and multi-domain learning with residual convolutional networks for emotion recognition. EB/OL.　2020-05-22. https://arxiv. org/pdf/1802. 06664. pdf

Pons G and Masip D. 2018b. Supervised committee of convolutional neural networks in automated facial expression analysis. IEEE Transactions on Affective Computing　9(3): 343-350　DOI: 10. 1109/TAFFC. 2017. 2753235

Reed S　Sohn K　Zhang Y and Lee H. 2014. Learning to disentangle factors of variation with manifold interaction//Proceedings of the 31st International Conference on International Conference on Machine Learning. New York: ACM: 1431-1439

Rifai S　Bengio Y　Courville A　Vincent P and Mirza M. 2012. Disentangling factors of variation for facial expression recognition//Proceedings of the 12th European Conference on Computer Vision. Florence: Springer: 808-822　DOI: 10. 1007/978-3-642-33783-3_58

Russell J A. 1980. A circumplex model of affect. Journal of Personality and Social Psychology　39(6): 1161-1178　DOI: 10. 1037/h0077714

Shan C F　Gong S G and McOwan P W. 2009. Facial expression recognition based on local binary patterns: a comprehensive study. Image and Vision Computing　27(6): 803-816　DOI: 10. 1016/j. imavis. 2008. 08. 005

Simonyan K and Zisserman A. 2014a. Two-stream convolutional networks for action recognition in videos//Proceedings of the 27th International Conference on Neural Information Processing Systems. Cambridge: ACM: 568-576

Simonyan K and Zisserman A. 2014b. Very deep convolutional networks for large-scale image recognition　EB/OL.　2020-05-22. https://arxiv. org/pdf/1409. 1556. pdf

Sun N　Li Q　Huan R Z　Liu J X and Han G. 2019. Deep spatial-tem-

poral feature fusion for facial expression recognition in static images. Pattern Recognition Letters 119: 49-61 DOI: 10.1016/j. patrec. 2017.10.022

Szegedy C Liu W Jia Y Q Sermanet P Reed S Anguelov D Erhan D Vanhoucke V and Rabinovich A. 2015. Going deeper with convolutions//Proceedings of 2015 IEEE Conference on Computer Vision and Pattern Recognition. Boston: IEEE: 1-9 DOI: 10.1109/CVPR.2015.7298594

Tian Y I Kanade T and Cohn J F. 2001. Recognizing action units for facial expression analysis. IEEE Transactions on Pattern Analysis and Machine Intelligence 23(2): 97-115 DOI: 10.1109/34.908962

Tang Y C. 2015. Deep learning using linear support vector machines. EB/OL. 2020-05-01. https://arxiv.org/pdf/1306.0239.pdf

Valstar M F Mehu M Jiang B H Pantic M and Scherer K. 2012. Meta-analysis of the first facial expression recognition challenge. IEEE Transactions on Systems Man and Cybernetics Part B (Cybernetics) 42(4): 966-979 DOI: 10.1109/TSMCB.2012.2200675

Vielzeuf V Pateux S and Jurie F. 2017. Temporal multimodal fusion for video emotion classification in the wild//Proceedings of the 19th ACM International Conference on Multimodal Interaction. New York: ACM: 569-576 DOI: 10.1145/3136755.3143011

Wang S F Zheng Z Q Yin S Yang J J and Ji Q. 2020. A novel dynamic model capturing spatial and temporal patterns for facial expression analysis. IEEE Transactions on Pattern Analysis and Machine Intelligence 42(9): 2082-2095 DOI: 10.1109/TPAMI.2019.2911937

Wang S J Li B J Liu Y J Yan W J Ou X Y Huang X H Xu F and Fu X L. 2018. Micro-expression recognition with small sample size by transferring long-term convolutional neural network. Neurocomputing 312: 251-262 DOI: 10.1016/j.neucom.2018.05.107

Wei X F Li H B Sun J and Chen L M. 2018. Unsupervised domain adaptation with regularized optimal transport for multimodal 2D+3D facial expression recognition//Proceedings of 13th IEEE International Conference on Automatic Face and Gesture Recognition. Xi'an: IEEE: 31-37 DOI: 10.1109/FG.2018.00015

Wu B F and Lin C H. 2018. Adaptive feature mapping for customizing deep learning based facial expression recognition model. IEEE Access 6: 12451-12461 DOI: 10.1109/ACCESS.2018.2805861

Yan J W Zheng W M Cui Z Tang C G Zhang T Zong Y and Sun N. 2016. Multi-clue fusion for emotion recognition in the wild//Proceedings of the 18th ACM International Conference on Multimodal Interaction. New York: ACM: 458-463 DOI: 10.1145/2993148.2997630

Yan W J Li X B Wang S J Zhao G Y Liu Y J Chen Y H and Fu X L. 2014. CASME II: an improved spontaneous micro-expression database and the baseline evaluation. PLoS One 9(1): e86041

DOI: 10.1371/journal.pone.0086041

Yan W J Wu Q Liu Y J Wang S J and Fu X L. 2013. CASME database: a dataset of spontaneous micro-expressions collected from neutralized faces//Proceedings of 10th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition. Shanghai: IEEE: 1-7 DOI: 10.1109/FG.2013.6553799

Yang H Y Ciftci U and Yin L J 2018a. Facial expression recognition by de-expression residue learning//Proceedings of 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition. Salt Lake City: IEEE: 2168-2177 DOI: 10.1109/CVPR.2018.00231

Yang H Y Zhang Z and Yin L J. 2018b. Identity-adaptive facial expression recognition through expression regeneration using conditional generative adversarial networks//Proceedings of 13th IEEE International Conference on Automatic Face and Gesture Recognition. Xi'an: IEEE: 294-301 DOI: 10.1109/FG.2018.00050

Yi D Lei Z Liao S C and Li S Z. 2014. Learning face representation from scratch. EB/OL. 2020-05-31. https://arxiv.org/pdf/1411.7923.pdf

Yu Z D and Zhang C. 2015. Image based static facial expression recognition with multiple deep network learning//Proceedings of 2015 ACM on International Conference on Multimodal Interaction. New York: ACM: 435-442 DOI: 10.1145/2818346.2830595

Zeng N Y Zhang H Song B Y Liu W B Li Y R and Dobaie A M. 2018. Facial expression recognition via learning deep sparse autoencoders. Neurocomputing 273: 643-649 DOI: 10.1016/j.neucom.2017.08.043

Zhang F F Zhang T Z Mao Q R and Xu C S. 2018a. Joint pose and expression modeling for facial expression recognition//Proceedings of 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition. Salt Lake City: IEEE: 3359-3368 DOI: 10.1109/CVPR.2018.00354

Zhang K H Huang Y Z Du Y and Wang L. 2017. Facial expression recognition based on deep evolutional spatial-temporal networks. IEEE Transactions on Image Processing 26(9): 4193-4203 DOI: 10.1109/TIP.2017.2689999

Zhang T Zheng W M Cui Z Zong Y Yan J W and Yan K Y. 2016. A deep neural network-driven feature learning method for multi-view facial expression recognition. IEEE Transactions on Multimedia 18(12): 2528-2536 DOI: 10.1109/TMM.2016.2598092

Zhang X Zhang L Wang X J and Shum H Y. 2012. Finding celebrities in billions of web images. IEEE Transactions on Multimedia 14(4): 995-1007 DOI: 10.1109/TMM.2012.2186121

Zhang Z P Luo P Loy C C and Tang X O. 2018b. From facial expression recognition to interpersonal relation prediction. International Journal of Computer Vision 126(5): 550-569 DOI: 10.1007/s11263-017-1055-1

Zhao G Y and Pietikainen M. 2007. Dynamic texture recognition using local binary patterns with an application to facial expressions. IEEE Transactions on Pattern Analysis and Machine Intelligence 29(6):

915-928　DOI: 10. 1109 /TPAMI. 2007. 1110

Zhao G Y　Huang X H　Taini M　Li S Z and Pietikäinen M. 2011. Facial expression recognition from near-infrared videos. Image and Vision Computing　29(9): 607–619　DOI: 10. 1016 /j. imavis. 2011. 07. 002

Zheng W M　Zong Y　Zhou X Y and Xin X M. 2018. Cross-domain color facial expression recognition using transductive transfer subspace learning. IEEE Transactions on Affective Computing　9(1): 21-37　DOI: 10. 1109 /TAFFC. 2016. 2563432

Zhi R C　Flierl M　Ruan Q Q and Kleijn W B. 2011. Graph-preserving sparse nonnegative matrix factorization with application to facial expression recognition. IEEE Transactions on Systems　Man　and Cybernetics　Part B ( Cybernetics)　41 (1): 38-52　DOI: 10. 1109 /TSMCB. 2010. 2044788

Zhong L　Liu Q S　Yang P　Liu B　Huang J Z and Metaxas D N. 2012. Learning active facial patches for expression analysis//Proceedings of 2012 IEEE Conference on Computer Vision and Pattern Recognition. Providence: IEEE: 2562-2569　DOI: 10. 1109 /CVPR. 2012. 6247974

1995

。

E-mail: ls1995@ bupt. edu. cn

、

、　　、　　、

。

E-mail: whdeng@ bupt. edu. cn