

On the Pitfalls of Learning with Limited Data: A Facial Expression Recognition Case Study

Miguel Rodríguez Santander^{a,*}, Juan Hernández Albarracín^{a,*}, Adín Ramírez Rivera^{a,**}

^a*Institute of Computing, University of Campinas, Campinas, SP, Brazil*

Abstract

Deep learning models need large amounts of data for training. In video recognition and classification, significant advances were achieved with the introduction of new large databases. However, the creation of large-databases for training is infeasible in several scenarios. Thus, existing or small collected databases are typically joined and amplified to train these models. Nevertheless, training neural networks on limited data is not straightforward and comes with a set of problems. In this paper, we explore the effects of stacking databases, model initialization, and data amplification techniques when training with limited data on deep learning models' performance. We focused on the problem of Facial Expression Recognition from videos. We performed an extensive study with four databases at a different complexity and nine deep-learning architectures for video classification. We found that (i) complex training sets translate better to more stable test sets when trained with transfer learning and synthetically generated data, but their performance yields a high variance; (ii) training with more detailed data translates to more stable performance on novel scenarios (albeit with lower performance); (iii) merging heterogeneous data is not a straightforward improvement, as the type of augmentation and initialization is crucial; (iv) classical data augmentation cannot fill the holes created by joining largely separated datasets; and (v) inductive biases help to bridge the gap when paired with synthetic data, but this data is not enough when working with standard initialization techniques.

Keywords: Learning with limited data, Limited data, Video classification

1. Introduction

The development of computational intelligence techniques grew in recent years due to the broad adoption of deep learning models and increasing data availability. The amount of data used to train deep learning models is crucial to avoid model under- or over-fitting. Nevertheless, several domains (or problems) do not have large amounts of data available due to acquisition costs. Hence, in these limited scenarios, deep learning solutions are incipient or, if they already exist, their results are limited due to the lack of labeled data.

Although standard techniques, such as dropout (N. Srivastava et al., 2014) or batch normalization (Szegedy, Vanhoucke, et al., 2016), are widely used to overcome the well-known over- and under-fitting problems, the silver bullet for every generalization problem is increasing the quantity of training data. However, when no big data is available from a single source, there are several ways to increase data. On the one hand, one can agglomerate several small datasets, i.e., where data from similar databases is compiled. Another option is to do transfer learning, cross-domain adaptation, or cross-database training. Data on similar data-rich domains are used to create robust models that are subsequently fine-tuned on the limited-data problem. And there is data augmentation, where new training examples are derived from the limited dataset through an augmentation function (commonly, another neural network is used).

It is important to point out that some of the work that explore limited-data setups are mainly surveys. Thus, their models' quantitative comparison is of compiling nature (Lu et al., 2020; Shorten and Khoshgoftaar,

*Code available at <https://gitlab.com/mipl/learning-with-limited-data>.

*Equal contribution

**Corresponding author

Email addresses: miguel.rodriguez@ieee.org (Miguel Rodríguez Santander), juan.albarracin@ic.unicamp.br (Juan Hernández Albarracín), adin@ic.unicamp (Adín Ramírez Rivera)

2019; X. Wang et al., 2020). Although other works perform dedicated experiments to evaluate the impact of their techniques, the authors’ focus is on one particular augmentation/initialization technique to improve the performance, and not on their combinations nor their contribution to the final result. Only two works conducted an empirical study in which they combine more than one technique: Chen et al. (2019) model complexity and cross-domain generalization on recognized few-shot learning models, and Brigato and Iocchi (2020) study the effects of data augmentation and model complexity. Hence, there is no focus on understanding the impact of limited data in the performance of models.

To the best of our knowledge, the overwhelming majority of studies on the impact of techniques facing small data focus on improving their process to increase performance. However, there is no systematic exploration of the effects of the data in the pipeline. Thus, the generalization effects of applying data augmentation techniques or transfer learning to limited-data setups are not well understood. Therefore, our work is the first to study transfer learning impacts, dataset-stacking, data augmentation, and their combinations for a visual task.

In this work, we present the shortcomings of training with limited data and analyze the effects of different inductive biases related to data augmentation and model initialization techniques. We focus on the Facial Expression Recognition (FER) problem in videos, as it is a domain that lacks data in the majority of the available datasets and has a variety of setups (e.g., posed vs. non-posed expressions, accessories, lighting, etc.) that make the datasets significantly heterogeneous among them. We experiment across four FER datasets with different levels of complexity and nine deep-learning architectures for video recognition. Our study includes comparisons on the effectiveness between classical vs. semantic data augmentation, transfer learning vs. random initialization, 2D- vs. 3D-CNN architectures, and low- vs. high-variance data used for dataset stacking.

Our main results are related to the characteristics of the used datasets. For instance, the complexity of the training sets influences the learned model’s performance, but they are correlated with the method used to augment the data from them. On the one hand, complex training sets are better with transfer learning and synthetic data generators at the cost of higher variance in the results. On the other hand, more straightforward datasets produce more stable results regardless of the augmentation used, but the performance is lower than the fine-tuned methods. Another main result is that merging or stacking datasets to create a large database is not necessarily an improvement. This result follows from the gaps that appear in the training data modes that challenge the models’ training. To avoid these holes in the stacked data, one needs to take special care to initialize and augment the data (which translates to better final performance). Moreover, we observed that classical data augmentation techniques fell short, avoiding the previously mentioned problems. However, inductive biases help to bridge this gap when paired with synthetic data.

2. Related Work

Previous work devoted to study the impact of common techniques to cope with small data is not scarce (Altan, Karasu, and Bekiros, 2019; Brigato and Iocchi, 2020; Chen et al., 2019; Guo et al., 2020; Lu et al., 2020; Pinetz et al., 2019; Shijie et al., 2017; Shin et al., 2016; Shorten and Khoshgoftaar, 2019; Soekhoe et al., 2016; Tajbakhsh et al., 2016; H. Wang et al., 2020; X. Wang et al., 2020), and focuses on understanding the impact of either Transfer Learning techniques (Chen et al., 2019; Lu et al., 2020; Shin et al., 2016; Soekhoe et al., 2016; Tajbakhsh et al., 2016), data augmentation (Altan and Karasu, 2020; Brigato and Iocchi, 2020; Lu et al., 2020; Pinetz et al., 2019; Shijie et al., 2017; Shorten and Khoshgoftaar, 2019; X. Wang et al., 2020), or cross-domain adaptation (Chen et al., 2019; Guo et al., 2020; Lu et al., 2020; H. Wang et al., 2020). Other works that proposed semantic data augmentation or transfer learning approaches punctually explore such an impact in a standard experimental setup: they compare their performance against baselines without the proposed data augmentation or transfer learning solution (Bowles et al., 2018; Bozorgtabar et al., 2019; Huynh-The and D.-S. Kim, 2019; Jena et al., 2020; Menon et al., 2019; Milicevic et al., 2018; Y. Wang et al., 2019; R. Zhang et al., 2018; X. Zhang et al., 2019; Y. Zhang et al., 2019). Naturally, such works are not extensive because they usually do not explore many datasets or various architectures. Nevertheless, these works’ primary focus is to advance their respective tasks instead of understanding the overall effect of the techniques and the limited-data. Thus, there is a need to understand the effect that training with limited data has on the overall learning framework and the results that limited-data brings.

In this section, we present existing solutions to cope with limited-data while training deep neural networks. The effects of these solutions are the object of study in this paper, namely, data augmentation and transfer learning. Then, we pose the problem of Facial Expression Recognition through the lens of limited data and present a brief state of the works that deal with the issue through such a lens.

2.1. Data Augmentation

Data augmentation consists in creating new training data through a set of transformations on the existing data. This technique has become one of the most popular for training deep architectures. In this work, we use image processing techniques and synthetic data generation for data augmentation.

2.1.1. Classical Data Augmentation

Classical Data Augmentation increases the training data using simple transformations. For images, such transformations include flipping, rotation, cropping, noise injection, and changes in lighting, among others (He et al., 2016; G. Huang et al., 2017; Simonyan and Zisserman, 2015; R. K. Srivastava et al., 2015). This technique is usually implemented online; i.e., once a training batch of the original data is sampled; each example has some probability of being transformed through one of the operations above. Nowadays, modern data augmentation methods automatically select existing techniques (Cubuk et al., 2019), corrupt features as augmentation (Maaten et al., 2013; Y. Wang et al., 2019), or apply class identity preserving transformations (Jaderberg et al., 2016; Ratner et al., 2017).

2.1.2. Synthetic Data Augmentation

Other works approached the lack-of-data problem using Generative Adversarial Networks (GANs) to perform semantic data augmentation on the training data, reducing the models' over-fitting. They applied it to a set of common problems, such as classification (Frid-Adar et al., 2018; Perez and J. Wang, 2017), segmentation (Bowles et al., 2018), detection (Han, Murao, et al., 2019; Han, Rundo, et al., 2020), among others. Beyond these tasks, training GANs with limited data is challenging since generating realistic results with small datasets becomes impossible due to the over-fitting problems in the neural network. To solve these problems, for instance, M. Zhao et al. (2020) used transfer learning techniques to train GANs in problems with limited data. On the other hand, Karras et al. (2020) proposed a novel technique to reduce the discriminator over-fitting through strategies that reduce noise when adding classic data augmentation processes.

The advances in Deep Generative Models, in particular GANs (Arjovsky et al., 2017; Brock et al., 2019; Goodfellow et al., 2014; Gulrajani et al., 2017), have made it possible to synthesize data with increasing realism. In the video synthesis domain, recent works in video reenactment have made it possible to generate videos whose objects (or characters) of interest mimic other videos (Bansal et al., 2018; Chan et al., 2019; Siarohin et al., 2019; L. Zhao et al., 2018). In particular, works on face reenactment (Aberman et al., 2019; Nirkin et al., 2019; Wu et al., 2018; Zakharov et al., 2019) have attained realistic results.

Although these methods have the potential for taking data augmentation for video to a higher semantic level, this idea has been widely explored mostly in images (Shorten and Khoshgoftaar, 2019; Y. Wang et al., 2019). Recent works started augmenting data based on GANs for video classification, with dynamical images, i.e., only one frame representing the whole video (S. Zhang et al., 2019). For this setting, we chose Monkey-Net by Siarohin et al. (2019) as the model to synthesize novel videos, due to its simplicity w.r.t. other reenactment methods and its self-supervised learning mechanism, i.e., the architecture learns keypoints and optical flow while learning to reenact videos.

2.2. Transfer Learning

Under the analogy of human cognition's capacity to transfer knowledge from one domain to another, Transfer Learning involves techniques to train models on tasks where large amounts of data are available and use the learned parameters in tasks where data is scarce. With the growth of deep learning, and the easiness of reusing models in different tasks, transferring techniques, (such as fine-tuning) have become the basis of modern Transfer Learning. Training models with large data sets allows the first layers to learn to extract generic features of the data, while the last layers extract specific features of the task to be solved. Fine-tuning was born from this idea, where already-trained models are used and retrained, either entirely or only a subset of its layers, in order to adapt to the new task model.

In this work, we use the fine-tuning technique by re-training all layers of the models that were already pre-trained with massive databases—e.g., ImageNet (J. Deng et al., 2009), UCF-101 (Soomro et al., 2012), Kinects-700 (Carreira and Zisserman, 2017). This technique serves as an inductive bias on the models since it assumes that the parameters learned for large-dataset image recognition are useful for the FER problem in videos from small datasets.

2.3. Facial Expression Recognition as a Limited Data Problem

Facial expression recognition (FER) has been widely studied in, what we pose as, a limited data environment (Y. Huang et al., 2019; Li and W. Deng, 2018). Existing approaches involve recognizing the facial expression through static images (Acharya et al., 2018; Kuo et al., 2018; Lopes et al., 2017; Marrero Fernandez et al., 2019; Ramírez Rivera, Rojas Castillo, et al., 2013; Ryu et al., 2017) or videos (Acharya et al., 2018; Daizong Liu, 2020; Hasani and Mahoor, 2017a; Hasani and Mahoor, 2017b; Jung et al., 2015; Kaya et al., 2017; Klaser et al., 2008; Kuo et al., 2018; Liang et al., 2019; M. Liu, Shan, et al., 2014; M. Liu, Li, et al., 2015; Ramírez Rivera and Chae, 2015; Ramírez Rivera, Rojas Castillo, et al., 2015; Yan, 2018; K. Zhang et al., 2017; S. Zhang et al., 2019; T. Zhang et al., 2019; G. Zhao, X. Huang, et al., 2011; G. Zhao and Pietikainen, 2007). The latter is incredibly difficult due to lack of data, since existing databases (Aifanti, Papachristou, et al., 2010; Dhall et al., 2011; Dhall et al., 2012; Kanade et al., 2000; Lucey et al., 2010; Pantic et al., 2005; Valstar and Pantic, 2010; G. Zhao, X. Huang, et al., 2011) have a limited amount of videos depicting the expressions. While image-based methods have an advantage due to a larger amount of frames extracted from these databases, video-based methods struggle with a restricted set of data from which to learn the patterns. On the other hand, methods that rely on handcrafted features (Kaya et al., 2017; B.-K. Kim et al., 2016; S. Zhang et al., 2019) need fewer data to learn in contrast to fully neural-network-based end-to-end learning methods (Acharya et al., 2018; Ding et al., 2017; Hasani and Mahoor, 2017a; Hasani and Mahoor, 2017b; Jung et al., 2015; Kuo et al., 2018; Liang et al., 2019; M. Liu, Li, et al., 2015; Marrero Fernandez et al., 2019; Mollahosseini et al., 2016; K. Zhang et al., 2017; X. Zhao et al., 2016). Thus, they tend to obtain better results in these limited data scenarios due to the inductive biases used on the descriptors.

The vast majority of the literature’s solutions focus on performance metrics on the test set regardless of their generalization to other domains. This metric chasing leads to methods that cannot be used in environments with characteristics not covered by the training database. One of the methodologies used to demonstrate such generalization capacity is the cross-database validation that consists in carrying out a proposed model’s training with one database and the validation on a different one.

The main problem with training deep learning models with small databases is that they tend to over-fit the training data and get poor results when performing tasks in real environments. In the last few years, many large databases have been released for public use, which allowed the development of new models that are more efficient and with low over-fitting rates. Models pre-trained with large databases are used to solve other literature problems that do not have large databases. Unlike natural-image databases (many of which contain millions of examples), there is a significant amount of databases of facial expressions in video available for public use. Still, all of them have a limited amount of data.

3. Experimental Design

We are interested in understanding the impact of the standard techniques to aggregate and augment data for small datasets, which generally are not enough to train deep models without the risk of over-fitting. One common approach is to create samples from the existing databases artificially (i.e., data augmentation). Another approach is to include inductive biases to the model, in the form of pre-trained weights (i.e., fine-tuning or transfer learning). Another approach consists of stacking databases from different setups to increase (i.e., augment) the data. These approaches can be either applied independently or combined in the same training pipeline.

Fine-tuning is commonly used in several methods in the literature. It is easy to start with existing models and further tune them for related problems. However, there is no attention paid to this transfer’s limits in the literature when performed on limited data. Moreover, novel methods with enough computing power (and data) rely on random initialization to avoid introducing biases from prior problems or data.

Data augmentation has been less studied but widely used in practice when no sustained effort to produce large amounts of data exist. Another practice is to pool several databases (or data capture setups) together.

However, this practice introduces heterogeneous biases and sources of error that are not leveraged by many samples that can be used reliably to train the model.

Our main objective is to analyze the effect of these two sources of error when training with limited data, i.e., how to initialize the model, and how to augment the data. In particular, we selected the facial expression recognition problem as it has several databases with a limited amount of samples and different setups and variable challenges. Thus, stacking these databases will introduce different levels of uncertainty and variability that the models may not handle.

We propose to contrast deep learning classifiers for video (based on 2D and 3D convolutions and with different recurrent mechanisms) by using two initialization methods, two data augmentation techniques, and stacking the databases in different ways (using a hold-out database and hold-out cross-validation folds to get other ideas of the generalization). Then, we discuss the observed results and conclude them.

3.1. Proposed Experiments

We executed a series of experiments to compare different network setups’ performance and their generalization when faced with limited data.

Model Comparison. The objective of this experiment is to obtain state-of-the-art models’ training metrics when training with limited data. We decided to perform an ablation between each model using different types of parameter initialization—Xavier (Glorot and Bengio, 2010) and transfer learning—and different types of data augmentation—classical and GAN-based. Each of these comparisons was made for each database used in the study.

Model Generalization. To obtain metrics on the generalization of our models, we decided to perform cross-database experiments. These experiments serve to compare model performances and assess their generalization capability. To measure the level of generalization, we proposed two cross-database experiments:

- Classic cross-database experiment: we trained with a database and tested on the others separately.
- All k -fold cross-validation: we created a k -fold cross-validation experiment by mixing all the databases into one. We carefully maintained the parallel folds from other k -fold cross-validation experiments to have paired experiments for comparison. With this test, we obtained insights on the generalization from the data sources when stacking them for learning versus when learning on constrained sets.

3.2. Data Sets

For this study, we consider four datasets whose size does not reach the order of thousands:

- **Extended CK+** (Lucey et al., 2010). It consist of 593 sequences with 123 subjects with seven emotion labels (anger, contempt, disgust, fear, happiness, sadness, and surprise) recorded on a controlled environment.
- **MMI** (Pantic et al., 2005; Valstar and Pantic, 2010). It contains 205 sequences from 30 subjects with six labels. This database is not recorded in the wild, but the subjects present variations of worn accessories and movements during the expression performance.
- **MUG** (Aifanti, Papachristou, et al., 2010). It contains 897 RGB videos of 52 different subjects and six expression classes: anger, disgust, fear, happiness, sadness, and surprise. All the subjects are performing all the expressions.
- **OULU-CASIA** (G. Zhao, X. Huang, et al., 2011). It comprises six expressions on three illumination conditions with 480 sequences per setting (80 subjects with six expressions each).

3.3. Classical and Deep Learning Architectures

For this work, we decided to compare (exhaustively, to the best of our efforts) the generalization capabilities of existing (and mostly used) deep learning architectures by training them with limited data. Since we are interested in video tasks, particularly for FER, we selected 2D convolutional networks paired with recurrent models, 3D convolutional networks that deal with the sequence as a whole, and 3D convolutional networks paired with recurrent models to handle the sequence.

Table 1. Hyperparameters of the commonly used architectures that we included in our study.

Model name	Type	Frame size	Batch size	Layers	LSTM size	GPUs	Parameters
VGG16-LSTM	Conv2D	(224,224)	10	16	1024	1	155M
InceptionV3-LSTM	Conv2D	(300,300)	3	47	512	2	33M
ResNet18-LSTM	Conv2D	(224,224)	10	18	512	1	13M
ResNet101-LSTM	Conv2D	(224,224)	3	101	512	2	46M
C3D	Conv3D	(100,100)	10	10	–	1	78M
I3D	Conv3D	(226,226)	5	47	–	2	12M
ResNet3D-18	Conv3D	(100,100)	10	18	–	1	33M
ResNet3D-101	Conv3D	(100,100)	10	101	–	1	85M
C3D-Block-LSTM	Conv3D	(100,100)	10	10	1024	1	66M

To assess the difference of performance between 3D-CNN-based architectures and LSTM-based ones, we opted for well-known “vanilla” architectures, with repeatable yet straightforward data setups. Our objective is to make comparable scenarios, unlike some state-of-the-art works that lack open code or have complex loss functions with incomplete information that difficult a homogeneous comparison due to different training setups.

3.3.1. Two-Dimensional Convolutional LSTM-based Architectures

The 2D convolutional architecture was designed for image classification and commonly trained on ImageNet (J. Deng et al., 2009). We decided to use VGG16 (Simonyan and Zisserman, 2015), Inception V3 (Szegedy, Ioffe, et al., 2016; Szegedy, W. Liu, et al., 2015; Szegedy, Vanhoucke, et al., 2016), and ResNet-18 and -101 (He et al., 2016) due to their widespread use. To adapt these architectures to work with videos, it was necessary to add recurrence. We opted to use an LSTM (Hochreiter and Schmidhuber, 1997) that receives the extracted feature vectors by the 2D convolutional network for each frame and then calculate the recurrence relationships between these vectors, obtaining a classification label for the video.

3.3.2. Three-Dimensional Convolutional Architectures

The 3D convolutional architectures were designed for video tasks, especially video classification. We decided to use C3D (Tran et al., 2015), I3D (Carreira and Zisserman, 2017), and ResNet3D-18 and -101 (Hara et al., 2018) since they obtain excellent results on UCF-101 (Soomro et al., 2012) or Kinetics (Carreira and Zisserman, 2017). These models work over the whole video, so there is no need to add a recurrence.

3.3.3. Three-Dimensional Convolutional LSTM-based Architectures

The main difference between the two types of architecture presented above is the use of recurrent cells. Three-dimensional convolution-based architectures extract spatiotemporal features with fixed time windows, while 2D-convolution-LSTM-based architectures extract spatial and temporal features independently, without time windows size limitations. With that in mind, we propose the use of smaller temporal windows for the 3D convolutions and further process the encodings with a recurrent cell. For this, we propose to use the C3D (Tran et al., 2015) architecture with an LSTM (Hochreiter and Schmidhuber, 1997) cell at the end. We call this architecture C3D-Block-LSTM.

3.4. Parameters

The hyperparameters used in each network are the same ones presented by the original authors. In models based on two-dimensional convolutions, an LSTM layer was added, followed by a classifier with three dense layers, where the last layer output is equal to the same number of classes from the database used. We show the model-dependent parameters in Table 1 (for implementation details of each model see the Appendix A). We trained them using Adam (Kingma and Ba, 2014) with a learning rate of 1×10^{-5} , weight decay of 5×10^{-3} , dropout (N. Srivastava et al., 2014) with probability of 0.8, and cross-entropy as a loss function. We used 5-fold cross-validation to obtain each model’s metrics by dividing each database into five different subgroups whose divisions were person-independent. We trimmed the videos using the most representative 25 frames of the expression except for the C3D and I3D architectures that used 16 and 64 frames, respectively (see Fig. 1 for the most representative parts on each database). We used PyTorch 1.0 (Paszke et al., 2017) to code and train all the models. The hardware used for the experiments were two NVIDIA Titan Xp graphic cards.

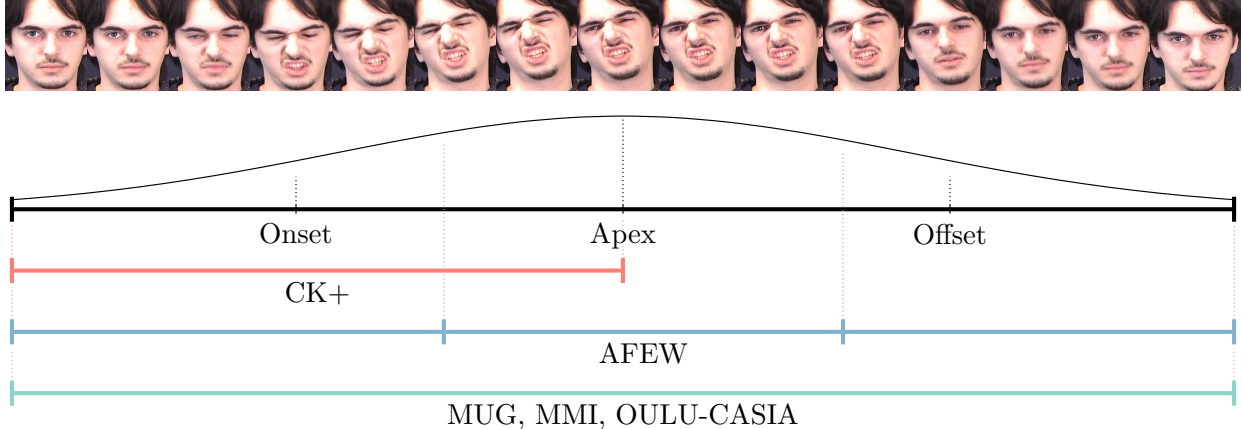


Figure 1. The evolution of facial expressions in videos commonly presents a Gaussian-bell-like behavior, in which each sequence starts with a neutral expression and then evolves (onset) to a determined facial expression (apex) and returns to neutral (offset). The observed temporal evolution is different in each database. E.g., sequences in CK+ (Kanade et al., 2000; Lucey et al., 2010) ends at the apex, while AFEW (Dhall et al., 2011; Dhall et al., 2012) contains clips showing different parts of the evolution, and MUG (Aifanti, Papachristou, et al., 2010), MMI (Pantic et al., 2005; Valstar and Pantic, 2010) and OULU-CASIA (G. Zhao, X. Huang, et al., 2011) show the full evolution.

4. Model Comparison with Limited Data

The experimental setup proposed for this study is based on three scenarios, which are: (i) single-database (Section 4.1), (ii) merged-database (Section 4.2), and (iii) cross-dataset (Section 4.3). For each scenario, we provide a four-dimensional analysis, in which each dimension corresponds to a family of variations on the training scheme. These dimensions are:

1. **Model Initialization.** As introduced above, we explore Random Initialization (RI) through Xavier (Glorot and Bengio, 2010), and Fine-tuning (FT).
2. **Data Augmentation Scheme.** We explore Classical Data Augmentation (DA), and Semantic Data Augmentation with Synthesized Data (SD).
3. **Models.** We use the nine models introduced above: VGG16-LSTM, InceptionV3-LSTM, ResNet18-LSTM, ResNet101-LSTM, C3D, C3D-Block-LSTM, I3D, ResNet3D-18, and ResNet3D-101.
4. **Databases.** As described above, we use four datasets: CK+, MMI, MUG, and OULU.

Due to the large volume of results obtained, we provide tables with detailed results in Appendixes B, C, and D, along with the results of exhaustive statistical tests in Appendix E. We compare the performances of different models setups in Figs. 2, 3, 4, and 5.

Every configuration in the setup described above (i.e. every initialization-augmentation-model-database combination) underwent a 5-fold cross-validation scheme. As suggested in the literature (Demšar, 2006; Dietterich, 1998), we performed Wilcoxon Signed-Rank tests (Demšar, 2006) between every pair of experiments, at 5% significance level. We show the p -values of those tests in Appendix E.

The lack of transparency regarding model configuration and training details from most of the literature methods prevents us to safely compare the performance of the models used in this study with the ones that conform to state of the art. Aspects like the train/validation/test split of the datasets, the number of folds used for cross-validation, the spatiotemporal cropping to pre-process the sequences, and the number of frames used for fixed-length models, must be aligned along with all the models for a fair comparison and few works report these details.

For that reason, a comparison between our performance metrics and the ones reported in the literature would be merely speculative and should not be considered as a proper comparison.

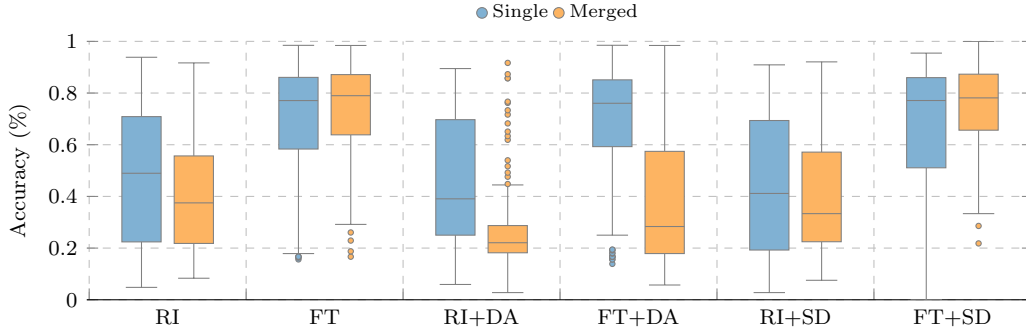


Figure 2. Aggregated results of Fig. 4 per training configuration. We compare training and evaluation on the same database (single) and using a merged database (merged). We used random (RI) and transfer learning by fine-tuning previously trained models (FT) for initialization. We used affine transformations used in classical data augmentation (DA) and synthetic data generated with models from the training partitions (SD) for augmentation.

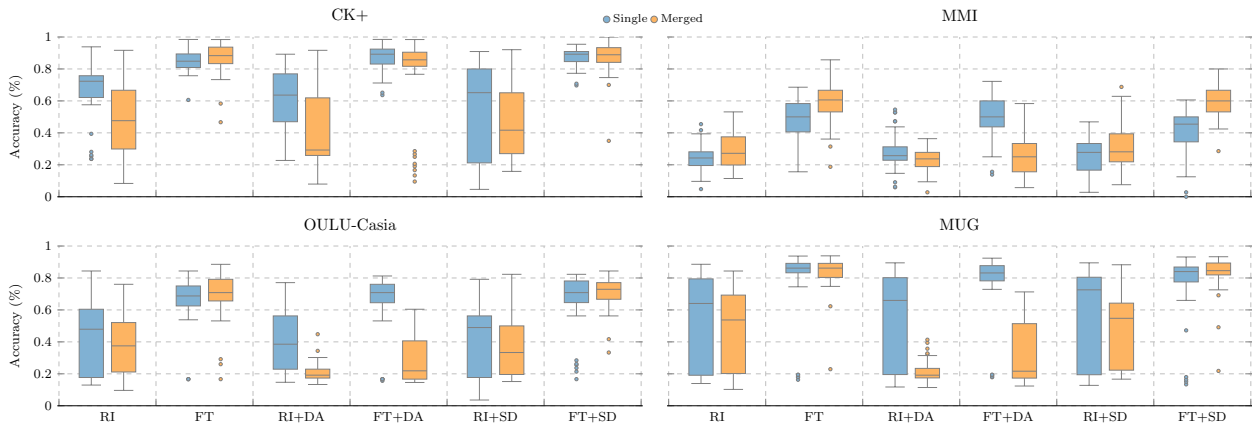


Figure 3. Aggregated results of Fig. 4 per dataset. We compare training and evaluation on the same database (single) and using a merged database (merged). We used random (RI) and transfer learning by fine-tuning previously trained models (FT) for initialization. We used affine transformations used in classical data augmentation (DA) and synthetic data generated with models from the training partitions (SD) for augmentation.

4.1. Single-Dataset k -Fold Cross-Validation

For the single-database scenario, we replicated standard experiments in the literature performed within a single database.

Figures 2 and 3 show aggregated results to compare experimental setups (i.e., initialization plus augmentation) for all and each dataset, respectively. The former aggregates paired experiments per method and database, while the latter aggregates only per method. The results correspond to the mean and variance of the 5-fold cross-validation setup. We can see that all configurations with FT obtained significantly better results than any configurations with RI. In general (see Fig. 2), the augmentation methods (either DA or SD) seems to be irrelevant in RI (cf. Table E.16). The only difference was observed in CK+, in which no augmentation presents a significant superiority over SD. For FT, MUG showed a decrease in performance when using any data augmentation, MMI showed that DA and no augmentation achieved better results than SD, and CK+ only showed the superiority of SD over no augmentation. OULU did not show any pattern. Diving into more detail, the left column of Fig. 4 shows the accuracy on each database by training the methods (shown with different colors) with several combinations of initialization-augmentation techniques (shown as groups within each database).

In summary, FT obtains better results than RI in all setups. For RI, the impact of any data augmentation is virtually null (except for CK+). On the other side, FT presents a more heterogeneous behavior w.r.t. the data augmentation technique and varies according to the dataset.

4.1.1. Model Initialization

As a first observation, the initialization method dramatically influences the results. When using RI, the models obtain lower and less stable metrics (higher variance) when compared to FT, which shows lower variances, in general. Compare the RI and FT groups in the left of Fig. 4 for each method, or aggregated in Figs. 2 and 3. Additionally, the temporal changes are a challenge to these methods. For instance, the CK+ and MUG databases are more stable because facial expressions are posed, unlike MMI and OULU. Furthermore, this result is reflected in the generalized higher classification performance in the former databases.

When using FT, the CNN-based networks reuse their learned features and tune them to the small amounts of data. It was evidenced that, for the 2D-CNN-based methods, FT yields significantly better performance than RI in all datasets, regardless of the type of data augmentation. On the other hand, the 3D-CNN-based methods present a more complex behavior since the improvement of FT w.r.t. RI seems to depend on the data augmentation type, and the results differ between datasets. For example, FT is significantly better than RI for CK+ and OULU when no augmentation is done and for DA, but this improvement is weaker for SD (cf. Fig. 2). FT improves in MMI when no data augmentation is used. For DA and SD, this improvement is weaker. Noticeably, the spread of variance is reduced and more consistent in MMI in contrast to CK+ or MUG. The latter presents the same behavior as MMI, except for the I3D, ResNet3D-101 and ResNet3D-18 models (cf. Fig. 4). I3D presented worst results for FT w.r.t. RI, and ResNet3D-101 and ResNet3D-18. These three methods presented a weaker improvement of FT for all the datasets, but only in MUG that was evident. Finally, note that by coupling the spatiotemporal features with a recurrence (i.e., a 3D block sequences processed through an LSTM on C3D-Block-LSTM), we obtained more unsatisfactory performance when compared with its purely 3D counterparts, and this difference is higher in the random initialization experiments.

When comparing 2D- versus 3D-based methods, we observed no significant difference in performance between the two families of methods. Individually, ResNet18-LSTM and C3D attained the lowest performances, while the rest of the methods were similar.

Analyzing the different models' results for their depth, we observe that models that have deeper architectures (i.e., VGG16-LSTM, InceptionV3-LSTM, I3D, ResNet101-LSTM, and R3D-101) present a subtle superiority w.r.t. the shallow architectures (i.e., C3D, C3D-Block-LSTM, ResNet18-LSTM, and ResNet3D-18). We observed that deep models took more advantage of fine-tuning than shallow ones, regardless of the data augmentation type, due to their bigger capacity, to which more data is beneficial.

In general, FT is more stable, i.e., with a smaller variance, in contrast to RI in the four datasets we tested (cf. Fig. 2). Moreover, FT obtains consistently better results than RI.

In summary, although FT presents a clear advantage w.r.t. RI, the latter presents a notably higher variance, and the significance of this improvement highly depends on the dataset, the model, and the augmentation technique. The most stable datasets, namely CK+ and MUG, obtained the best classification accuracies. For the 2D-CNN-based models, FT significantly outperforms RI in all datasets, regardless of the augmentation technique. On the other side, for the 3D-based models, this improvement is more heterogeneous: when no augmentation is performed, the superiority of FT over RI is significant but weaker for DA and SD. In general, there is no significant difference in the performance of 2D-CNN-based methods and 3D-CNN-based ones, while the LSTM+3D-CNN-method had a lower performance than the rest. Finally, deep architectures present a slightly superior performance over shallow ones for FT, while this difference is less evident for RI.

4.1.2. Data Augmentation

As described above, we considered two types of data augmentation: classical (DA) and GAN-based (SD). For DA, we obtained the former with affine transformations on the images (e.g., rotations, and flip), commonly used in previous work. For SD, we trained one generative model per experiment, namely Monkey-Net (Siarohin et al., 2019), from the training partition, and then used it to generate novel examples as a source for data augmentation. We selected these two augmentation methods due to their wide-spread use and limited ourselves to them due to computational resources to execute the experiments. Table B.1 shows the detailed comparison between the impact of DA and SD on the model performance.

The impact of DA for RI depends on both the database and the model. For example, for CK+, DA's effect in the C3D-block-LSTM, InceptionV3-LSTM, ResNet101-LSTM, ResNet3D-101, and VGG16-LSTM architectures is a decrease on accuracy performance that can be considered as significant, while the converse is true for C3D and I3D. For ResNet18-LSTM, no significant difference was evidenced (cf. Fig. 4). For OULU,

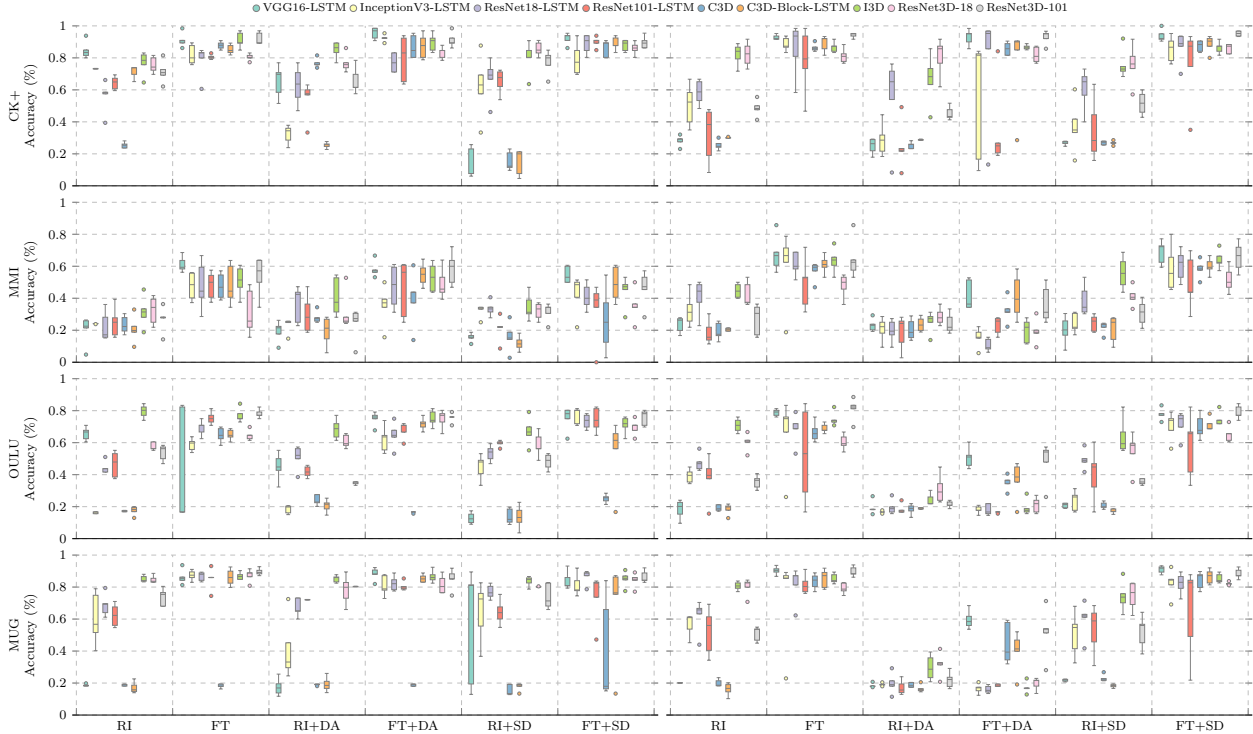


Figure 4. Accuracy of several models when trained with different combinations of initialization and augmentation techniques. We used random (RI) and transfer learning by fine-tuning previously trained models (FT) for initialization. We used affine transformations used in classical data augmentation (DA) and synthetic data generated with models from the training partitions (SD) for augmentation. We show training and evaluation on the same database (left) and using a merged database (right)—cf. Sections 4.1 and 4.2, respectively.

DA increased accuracy on C3D and decreased it for I3D, ResNet3D-101, and VGG16-LSTM. In general, for no database, RI presented any significant difference in performance when comparing RI+DA. Nevertheless, we observed an increment in the variance w.r.t. the non-augmented version in all datasets (cf. Fig. 2).

For FT, DA shows a non-significant general improvement (cf. Table E.13). Note, however, that for MUG, DA significantly reduced the accuracy of FT. There are more outliers in OULU (cf. Fig. 4 FT vs. FT+DA), on which C3D-block-LSTM and ResNet3D-18 got their accuracy increased, while C3D and ResNet101-LSTM got it decreased.

Regarding the impact of SD for RI, CK+ presented an increase in accuracy only for ResNet3D-18 and decreased for C3D, C3D-block-LSTM, and VGG16-LSTM. We observed no significant changes in the rest of the models. Similarly, we observed no significant impact for the MUG dataset, except for a decrease in the performance of ResNet3D-18. OULU presented increments in InceptionV3-LSTM, ResNet101-LSTM and ResNet18-LSTM, and decrements in I3D and VGG16-LSTM. MMI presented an increment in InceptionV3-LSTM and a decrement in C3D-block-LSTM. In general, for each dataset, there is no significant change when synthetically augmenting the RI.

For FT, SD impacts the performance of the models less. CK+ only presented a significant increment in ResNet101-LSTM, ResNet18-LSTM, and ResNet3D-18. For MMI, VGG16-LSTM decreased its accuracy, and no other method presented significant differences. For MUG, ResNet3D-18 decreased its accuracy, and no other method presented significant differences. OULU presented and incremented in InceptionV3-LSTM and a decrement in C3D. In general, in MMI, the SD also decreased its accuracy significantly compared to DA and without augmentation. And in MUG, no augmentation performed significantly better than the other two (cf. Table E.13).

When comparing DA vs. SD in RI, MMI did not present any significant difference while, for CK+, C3D, C3D-block-LSTM, and VGG16-LSTM presented a significant superiority of DA over SD. For MUG, InceptionV3-LSTM and ResNet18-LSTM presented a superiority of SD over DA. For OULU, C3D and

Table 2. The literature methods’ average accuracy and selected methods on our experiments on 5-fold cross-validation on single- and merged-databases. We highlight in bold the best results of each experiment set.

Model	CK+	MMI	OULU	MUG
Liang et al. (Liang et al., 2019)	99.6	80.7	91.0	–
Acharya et al. (Acharya et al., 2018)	–	–	–	–
Kuo et al. (Kuo et al., 2018)	98.4	–	91.6	–
Yan et al. (Yan, 2018)	96.6	–	–	–
Ghimire et al. (Ghimire et al., 2017)	97.8	77.2	–	95.5
Inception-ResNet-3D (Hasani and Mahoor, 2017a)	93.2	77.5	–	–
Kaya et al. (Kaya et al., 2017)	98.3	70.3	–	–
MSCNN-PHRNN (K. Zhang et al., 2017)	98.5	81.1	86.2	–
DTAGN (Jung et al., 2015)	97.2	70.2	81.4	–
Aifanti et al. (Aifanti and Delopoulos, 2014)	94.3	–	–	92.8
3DCNN-DAP (M. Liu, Shan, et al., 2014)	92.4	63.4	–	–
Single-Dataset 5-Fold Cross-Validation				
VGG16-LSTM-FT	91.1	61.2	43.1	86.2
VGG16-LSTM-FT-DA	95.4	58.2	75.0	88.2
I3D-RI	76.7	31.4	79.5	85.1
I3D-FT-DS	87.8	44.5	71.0	85.0
ResNet3D-101-FT	92.0	53.8	78.1	89.5
Merged-Dataset 5-Fold Cross-Validation				
ResNet3D-101-FT	94.2	64.6	81.0	89.7
ResNet3D-101-FT-SD	95.1	65.9	79.0	88.4
VGG16-LSTM-FT-SD	93.8	68.7	77.9	90.9

VGG16-LSTM presented a superiority of DA over SD, while InceptionV3-LSTM, ResNet101-LSTM, and ResNet3D-101 presented the opposite.

In FT, only a few significant differences were noticed. MUG did not present any significant improvement while, for CK+, ResNet18-LSTM showed significant superiority of SD over DA, and VGG16-LSTM presented the opposite. For MMI, ResNet3D-18 presented a significant superiority of DA over SD. For OULU, C3D presented a superiority of SD over DA, while C3D-block-LSTM presented the opposite.

Note that the classical data augmentation techniques (i.e., affine transformations) obtained the best result when selecting the maximum among the methods in most databases. However, in general, with a 5% significance level the FT methods outperformed them (cf. Table E.16). Since the affine transformations explore different local features of the original data, it generalizes well to the changes in the current database, i.e., it creates an invariant model to the changes within the subjects in the database. In contrast, the synthetic data generation explores new combinations within the database that generalize better, for instance, when trained and tested with diverse data (cf. Section 4.2) or when tested on different data to evaluate the method’s generalization capabilities (cf. Section 4.3). Due to the lack of challenging data within a single database, this benefit is not reflected in this experiment, but it becomes evident in the following experiments. CK+ and OULU, respectively, the simplest and the most complex databases, presented significant impacts of using the data augmentation techniques. Curiously, such impact was more evidenced in shallow models for CK+ and deep models for OULU. We will resume this discussion in the following sections.

4.1.3. Comparison Against Existing Methods

In Table 2, we present our highest results and existing experiments in the literature. We show the detailed results on Table B.1 in Appendix B. We highlight that our experiments are not directly comparable with those methods. Even our single-database k -fold cross-validation, which is the closest to the reported methods, do not use the reported hyper-parameters for all the methods (see Table 1). Recall that we synchronized the hyper-parameters and setups as much as possible to allow paired experiments. Hence, we cannot consider our results a fair nor a direct comparison against the literature.¹ Nevertheless, we show the literature results as an upper bound to our metrics and contrast them.

As shown in Table 2, there is no consensus on the highest accuracy over the databases. In general, it seems that initializing weights and fine-tuning them improves the results and that 2D based methods have the

¹Note that it is common to see comparisons on the literature where the contrasted methods have a different set of parameters. Nevertheless, these mistakes are rarely discussed.

edge over the 3D ones for the CK+ and MMI databases. While for OULU and MUG, the 3D based methods show the highest accuracy. However, as discussed before in Section 4.1, these methods present high variability, drops in accuracy, and inconsistent results across databases. Thus, a maximum value may be a misleading indicator.²

In contrast to the literature, our best results are significantly below the maximum reported values (Ghimire et al., 2017; Kaya et al., 2017; Liang et al., 2019). Liang et al.’s (2019) method uses a two-stream method to process the spatial and temporal information of the videos, with bidirectional recurrence. Interestingly, they use classical data augmentation to enhance their results. The work of Ghimire et al. (2017) relies on geometrical features with support vector machines and boosting to learn robust classifiers. Moreover, Kaya et al.’s (2017) work uses audiovisual information to produce deep and dense visual features with audio features fused with a kernel extreme learning machine. However, these complex and ad-hoc networks were outside of the scope of our exploration. Nevertheless, we find it interesting to show the wide range of results and lack of consistency in the existing approaches that lead to exhaustive searches to find local maxima that may not generalize well (as discussed in Section 4.3).

4.2. Merged-Dataset k -Fold Cross-Validation

Another common approach to tackle limited data is to join existing databases together. In this section, we analyze the effects of merging databases to increase the amount of data.

We created a new database by combining CK+, MMI, OULU-Casia, and MUG. We created a k -fold cross-validation experiment where each i th fold is the union of the i th folds at each database when doing the k -fold cross-validation independently. We took care of using the same folds of previous experiments to compare between merged and individual databases. I.e., each model-initialization–augmentation combination was trained once for each merged-fold of the 5-fold cross-validation. We obtained the test mean accuracy for each database so that each boxplot model-initialization–augmentation is comparable against its single-database counterpart.

In Fig. 3, we observe the same general dominance of FT over RI in the single-database experiments. MUG and OULU showed a significant dominance of RI with no augmentation and SD over FT with DA. When using FT, all the datasets showed a significant dominance of no augmentation and SD over DA. Similarly, when using RI, no augmentation dominates DA in all datasets and, for MMI, MUG and OULU, SD dominates DA.

In the right column of Fig. 4, we show the fine-grained average accuracy over the 5-fold cross-validation per database (remember that the training was done over a merged-database). We show the detailed experiments in Table C.1.

4.2.1. Model Initialization

Models such as C3D and VGG16-LSTM showed significant superiority of FT over RI, regardless of the data augmentation scheme. C3D-block-LSTM and ResNet3D-101 evidenced this superiority for SD and no augmentation at all. I3D and ResNet3D-18 experimented with no significant improvement of one initialization method w.r.t. the other, for any data augmentation scheme. Besides these two methods, no significant difference when using DA was evidenced for I3D, InceptionV3-LSTM, and ResNet101-LSTM. Only three cases reported superiority of RI over FT, all of them using DA: I3D on MUG, ResNet18-LSTM on MMI, and ResNet3D-18 on MUG.

In summary, the per-database improvement of FT over RI is significant (cf. Table E.14). And the former showed a decreased variance (cf. Fig. 3). And as a whole, the same trend was maintained (cf. Table E.17). However, when looked at a more fine-grained scale, we can see that the data augmentation technique dramatically affects the dominance of one initialization method over the other, and such behavior varies across the databases.

4.2.2. Data Augmentation

The experiments show a clear pattern related to the data augmentation method: DA decreases the models’ accuracy, regardless of the initialization (cf. Fig. 2). Furthermore, SD increases the variance of the results.

²A commonly reported indicator of the literature.

The apparent increment in the accuracy of the models when using SD cannot be considered as significant (cf. Table E.17). When comparing DA against SD, the latter trivially attains a better performance.

We observed only three cases of improvement of accuracy, and all were only for SD. VGG16-LSTM had an increment for RI and SD for MUG, and C3D-block-LSTM and ResNet3D-18 had an increment for FT for OULU. The two models whose accuracy remained invariant the most when applying DA with RI are C3D and VGG16-LSTM. However, this variance does not stand for FT, in which all the models significantly decreased their accuracy with DA.

It is interesting to notice that aggregating existing datasets—a somehow naive approach of augmenting data—may yield adverse outcomes due to the mixture and heterogeneity it induces. Additional gaps in data space may be too much to be filled by the models while learning, as demonstrated by the experiments. Augmentation techniques that traverse these data clusters may not be enough to close the gaps and instead hinder the overall accuracy. In contrast, synthetic generation helps. However, it is not a general solution since typical initialization yields lower results than the inductive bias in fine-tuned methods. Moreover, notice that the aggregated results (see Table E.17) cannot reject the hypothesis that each method and its synthetically augmented counterpart are equal.

4.2.3. Comparison Against Single-database Training

This section compares the classification performance of the single-database training against the merged-database training on the same test sets.

Figure 3 shows that the use of DA and FT implies a significant drop in performance for the merged dataset w.r.t. the single one, regardless of the initialization. RI also presents a drop in performance for the merged dataset, for CK+, MUG, and OULU, when using DA. It can be seen that, in most cases, the effect of the merged database is detrimental to the performance. However, if we compare RI on a single-database training with FT in the merged-database training and FT on a single-database training with RI in the merged-database training, it is evident that the configuration of FT attains the best results in both cases. Another remarkable pattern is that, for FT, SD consistently attains better results than DA. For RI, no pattern of this kind is observed.

In general (cf. Table E.18), all the methods in the single database reject the hypothesis that any RI method in the merge database is greater than them. In contrast, the FT methods on the merged database cannot reject the hypothesis that they are more significant than any single database method. The exception to the trend is FD+DA in the merged database that fails to reject the hypothesis that it is greater than the augmented versions of RI in the single database. Moreover, we notice an increase in the variance of the results in the different databases for at least one-third of the setups (cf. Fig. 2). In the merged database, particular exceptions for FT are C3D for MUG, ResNet3D-18 for MUG, and I3D for OULU. In the first, the merged database increases performance; while it decreases performance in the other two cases.

We observed a different behavior for FT when using classical data augmentation. For MUG and OULU, the merged database decreases the accuracy for all the models (cf. Fig. 2), except for C3D, in which it significantly increases. We observed no significant difference in CK+, except for InceptionV3-LSTM, in which the merged database decreased the performance. Regarding MMI, the majority of models experimented with a decrease in accuracy, except for C3D, C3D-block-LSTM, and InceptionV3-LSTM, which presented no significant difference. Notice that there is a large disparity between single and merged database results in the less constrained datasets. In contrast, CK+ contains more homogeneous images that have a smaller disparity in the results. Again, this reinforces the idea that heterogeneity in the data diminishes the gains when augmenting limited data problems.

Some models showed increments in their accuracy for FT with SD, such as C3D, C3D-block-LSTM, and ResNet3D-101. MMI showed increments for C3D, I3D, ResNet18-LSTM, ResNet3D-101 and VGG16-LSTM. In the rest of the cases, no significant difference was evidenced. In the general case, the FT+SD for the merged database, we can reject the hypothesis that it is smaller or equal to other methods.

When using random initialization, although the merged dataset’s effect depends on the model and the dataset, most cases showed either a decrease or no significant difference. Interestingly, random initializations do not take advantage of the augment in data (albeit with an increase in variance on the data), while the pre-trained models do.

As a remarkable pattern, we can say that MUG and OULU decreased their performance in general when using classical data augmentation, regardless of the initialization algorithm and the model used.

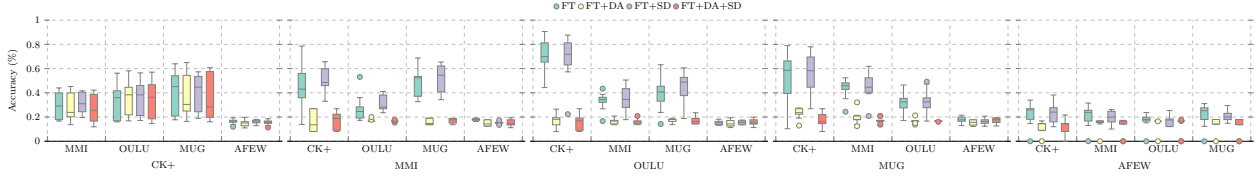


Figure 5. The average accuracy of all the architectures on a cross-database evaluation. Each plot represents the training database, and the different groups are the accuracy when evaluated on that database. For initialization, we used transfer learning by fine-tuning previously trained models (FT). We used affine transformations used in classical data augmentation (DA) and synthetic data generated with models from the training partitions (SD) for augmentation.

4.3. Cross-Dataset Cross-Validation

In this experiment, we did a cross-database evaluation, consisting of training on all the videos from one database and testing all the videos from all others. This experiment’s objective is to evaluate the generalization capabilities of the models when trained and tested on different domains.

Due to time and computational resources constraints, and based on the results obtained in the experiments presented previously, we decided to limit our cross-database experiments only to FT initialization since RI performed more poorly in all setups. For this set of experiments, we incorporated a new variation in combining both DA and SD augmentation techniques. We denote them as FT+DA+SD.

Additionally, for this generalization scenario, we decided to include the Acted Facial Expressions in the Wild (AFEW) database (Dhall et al., 2011; Dhall et al., 2012), which contains in-the-wild videos extracted from movies with seven spontaneous expressions. Unlike the previous databases that have no partition for train, validation, or test, AFEW is divided into three groups of 773, 383, and 653 samples for, respectively, train, validation, and test. This database is the most complex one, as it is the most heterogeneous.

In our generalization test, we found that DA yielded lower accuracy w.r.t. SD. See Fig. 5 for the database hold-out experiments. We show the detailed results in Appendix D, where each table represents a model-initialization-augmentation combination. Moreover, the addition of affine transformations to the synthetic data hurts the generalization capabilities of the models (cf. FT+DA+SD in Fig. 5).

Interestingly, when training on a stable and straightforward database, like CK+, the methods’ generalization is stable for the different augmentation techniques. On the contrary, when training in more diverse databases, like MMI, OULU, and MUG, the affine data augmentation does not work. In particular, when involving the AFEW database, either for training or testing, the attained results significantly dropped, regardless of the initialization and data augmentation strategy. Note that even though using SD is minor when compared with FT alone, the variance is reduced w.r.t. the former.

5. Discussion

In this section, we wrap up the comparisons of this study more conclusively by discussing the impact of each of the studied techniques.

5.1. Impact of Initialization Techniques

In this study, the most significant pattern observed is the superiority of inductive biases in the form of pre-trained models (a.k.a. FT) over random initializations (RI). The latter also presents higher variability in most of the cases for single-database training (Section 4.1). We observed the dominance of FT over RI to happen regardless of both the datasets and the data augmentation techniques. These two factors only affect the magnitude of such dominance. Our experiments showed no impact on database stacking over the choice of initialization. Furthermore, it showed clear superiority of FT with and without the dataset stacking. Consequently, one may prefer to opt for fine-tuning instead of stacking databases when facing small data (Section 4.2). We also observed that the improvement provided by FT is more significant in deeper architectures when tested on a single database (Section 4.1). RI also shows some detrimental effect on database stacking (Section 4.2).

5.2. Impact of Data Augmentation Techniques

The type of data augmentation (either DA or SD) plays a vital role in the methods’ performance. We saw that it impacts the effect of database stacking, although to a lesser extent than initialization, since the effectiveness of the augmentation technique is tightly constrained by the initialization (Section 4.2). As an example, notice that, for RI, the effects of both augmentation techniques, DA and SD, are practically null, while, for FT, they are extremely heterogeneous (Section 4.1). When training with a single database, the effects of DA are null for most of the datasets, but mainly on high-variance datasets like OULU and MMI. These results suggest that merely affine transformations do not augment the diversity of such datasets. When training on merged databases, DA generally decreases the performance of the models, while SD, besides not providing any significant improvement, increases the variance (Section 4.2). The only scenario in which SD significantly outperforms DA and no augmentation corresponds to the cross-database experiments. This result shows that synthetic data augmentation is an effective technique to increase the models’ generalization performance (Section 4.3).

5.3. Impact of Database Stacking

Augmenting data by stacking datasets is not guaranteed to be an improvement (Section 4.2). The results we obtained in this study show that naively mixing databases without a care for their contents and similarities may produce sparse training sets that tax the model to the point to decrease its performance. Simply said, if not carefully selected, the augmented data may increase the gaps of the base dataset instead of filling them. We assume that the decrement in performance is related to the models being incapable of not bridging the training data gaps. Notice that these gaps are not easily controllable without a care in understanding the data available. Moreover, we suggest working with limited data and better inductive biases instead of blindly aggregating data. And the final aggregated results (cf. Fig. 2) support that fine-tuned methods outperform their counterparts.

5.4. Impacts on Generalization through Cross-database Testing

Our generalization studies by cross-database testing show that data augmentation through affine transformations works better within the same data setup but decay when testing in heterogeneous data setups (Section 4.3). That is, the affine transformation has limited generalization capabilities. This behavior is explained by the fact that this technique only has affine transformations, i.e., it is only creating the same data presented in different shapes, which does not help the fit of the general sample space of real-world tasks. This limitation results in a database over-fitting, responsible for the difference in metrics between training with only one database and training using several databases.

We found that using stable data (with fewer variations) produces more predictable and stable results over newer data regarding cross-database evaluation. However, training in heterogeneous data, synthetic data, and fine-tuning helped improve performance in other scenarios, at the cost of a wider variability on the predicted performance.

6. Conclusions

We presented a large-scale study to show the improvements and limitations when training with limited data for classification problems dealing with facial expression recognition from video. We explored the use of two model initialization techniques and two data augmentation methods, along with the possibility of stacking databases. These variants were tested on nine widely-used deep architectures and in four datasets. We performed an exhaustive analysis of all these variants and saw that the improvement in classification performance that these techniques provide is not straightforward, but full of nuances. We showed how mixing these techniques yield different generalization levels and how significant the differences between them are in terms of classification performance.

Among the most overwhelming results obtained from the study, the dominance of inductive biases in the form of Fine-Tuning over other model initialization and data augmentation techniques stands out the most, suggesting that Fine-Tuning is the most likely technique to improve model performance when facing small datasets significantly. A second insight from our study is that data augmentation, either classical or GAN-based, marginally improves most of the tested models’ performance when no other technique (e.g.

database stacking or fine-tuning) is used. When database stacking is used, GAN-based data augmentation is the most suited to fill the gaps created by merging heterogeneous data, while classical data augmentation falls short. A third insight is that stacking heterogeneous databases is not a straightforward improvement, as the type of augmentation and initialization is crucial. Besides that, the complexity of the test data and the complexity of the stacking datasets to train a model are essential aspects we must consider before merging data indiscriminately, as complex training sets translate better to more stable test sets when fine-tuned and coupled with GAN-based-synthetic data. On the other hand, more simple data translates to more stable performance on unseen data (albeit with lower performance).

Acknowledgments

This project was funded by the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior—Brasil (CAPES)—Finance Code 001; São Paulo Research Foundation (FAPESP) under grants No. 2016/19947-6, 2017/16144-2, and 2019/07257-3; and by the Brazilian National Council for Scientific and Technological Development (CNPq) under grant No. 307425/2017-7.

References

- Aberman, K., Wu, R., Lischinski, D., Chen, B., & Cohen-Or, D. (2019). Learning character-agnostic motion for motion retargeting in 2d. *ACM Trans. Graphics*, 38(4), 1–14.
- Acharya, D., Huang, Z., Paudel, D. P., & Van Gool, L. (2018). Covariance pooling for facial expression recognition. *IEEE Inter. Conf. Comput. Vis., Pattern Recog. Wksp. (CVPRW)*, 480–4807.
- Aifanti, N., Papachristou, C., & Delopoulos, A. (2010). The MUG facial expression database. *Inter. Wksp. Image Anal. Multimedia Interact. Serv. (WIAMIS)*, 1–4.
- Aifanti, N., & Delopoulos, A. (2014). Linear subspaces for facial expression recognition. *Signal Process. Img. Comm.*, 29(1), 177–188.
- Altan, A., & Karasu, S. (2020). Recognition of COVID-19 disease from x-ray images by hybrid model consisting of 2d curvelet transform, chaotic salp swarm algorithm and deep learning technique. *Chaos, Solitons & Fractals*, 140, 110071.
- Altan, A., Karasu, S., & Bekiros, S. (2019). Digital currency forecasting with chaotic meta-heuristic bio-inspired signal processing techniques. *Chaos, Solitons & Fractals*, 126, 325–336.
- Arjovsky, M., Chintala, S., & Bottou, L. (2017). Wasserstein generative adversarial networks. *Inter. Conf. Mach. Learn. (ICML)*.
- Bansal, A., Ma, S., Ramanan, D., & Sheikh, Y. (2018). Recycle-GAN: Unsupervised video retargeting. *European Conf. Comput. Vis. (ECCV)*.
- Bowles, C., Chen, L., Guerrero, R., Bentley, P., Gunn, R., Hammers, A., Dickie, D. A., Hernández, M. V., Wardlaw, J., & Rueckert, D. (2018). GAN augmentation: Augmenting training data using generative adversarial networks. *arXiv preprint arXiv:1810.10863*.
- Bozorgtabar, B., Rad, M. S., Ekenel, H. K., & Thiran, J.-P. (2019). Using photorealistic face synthesis and domain adaptation to improve facial expression analysis. *IEEE Inter. Conf. Automat. Face Gesture Recog.*
- Brigato, L., & Iocchi, L. (2020). A close look at deep learning with small data. *IEEE Inter. Conf. Pattern Recog. (ICPR)*.
- Brock, A., Donahue, J., & Simonyan, K. (2019). Large scale GAN training for high fidelity natural image synthesis. *Inter. Conf. Learn. Represent. (ICLR)*.
- Carreira, J., & Zisserman, A. (2017). Quo vadis, action recognition? a new model and the kinetics dataset. *IEEE/CVF Inter. Conf. Comput. Vis. Pattern Recog. (CVPR)*, 4724–4733.
- Chan, C., Ginosar, S., Zhou, T., & Efros, A. A. (2019). Everybody dance now. *IEEE Inter. Conf. Comput. Vis. (ICCV)*.
- Chen, W.-Y., Liu, Y.-C., Kira, Z., Wang, Y.-C. F., & Huang, J.-B. (2019). A closer look at few-shot classification. *Inter. Conf. Learn. Represent. (ICLR)*.
- Cubuk, E. D., Zoph, B., Mane, D., Vasudevan, V., & Le, Q. V. (2019). Autoaugment: Learning augmentation strategies from data. *IEEE/CVF Inter. Conf. Comput. Vis. Pattern Recog. (CVPR)*.
- Daizong Liu, P. Z., Hongting Zhang. (2020). Video-based facial expression recognition using graph convolutional networks. *IEEE Inter. Conf. Pattern Recog. (ICPR)*.
- Demšar, J. (2006). Statistical comparisons of classifiers over multiple data sets. *J. Mach. Learn. Res.*, 7(Jan), 1–30.
- Deng, J., Dong, W., Socher, R., Li, L., Li, K., & Li, F.-F. (2009). ImageNet: A large-scale hierarchical image database. *IEEE/CVF Inter. Conf. Comput. Vis. Pattern Recog. (CVPR)*, 248–255.

- Dhall, A., Goecke, R., Lucey, S., & Gedeon, T. (2011). Acted facial expressions in the wild database. *Australian National University, Canberra, Australia, Technical Report TR-CS-11, 2*, 1.
- Dhall, A., Goecke, R., Lucey, S., & Gedeon, T. (2012). A semi-automatic method for collecting richly labelled large facial expression databases from movies. *IEEE Multimedia*.
- Dietterich, T. G. (1998). Approximate statistical tests for comparing supervised classification learning algorithms. *Neural Comput.*, *10*(7), 1895–1923.
- Ding, H., Zhou, S. K., & Chellappa, R. (2017). FaceNet2ExpNet: Regularizing a deep face recognition net for expression recognition. *IEEE Inter. Conf. Automat. Face Gesture Recog.*, 118–126.
- Frid-Adar, M., Diamant, I., Klang, E., Amitai, M., Goldberger, J., & Greenspan, H. (2018). GAN-based synthetic medical image augmentation for increased CNN performance in liver lesion classification. *Neurocomputing*, *321*, 321–331.
- Ghimire, D., Lee, J., Li, Z.-N., & Jeong, S. (2017). Recognition of facial expressions based on salient geometric features and support vector machines. *Multimedia Tools Appl.*, *76*(6), 7921–7946.
- Glorot, X., & Bengio, Y. (2010). Understanding the difficulty of training deep feedforward neural networks. *Inter. Conf. Artif. Intell. Stat. (AISTATS)*, 249–256.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., & Bengio, Y. (2014). Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, & K. Q. Weinberger (Eds.), *Adv. neural inf. process. sys. (NeurIPS)* (pp. 2672–2680). Curran Associates, Inc.
- Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., & Courville, A. C. (2017). Improved training of Wasserstein GANs. *Neural Inf. Process. Sys. (NeurIPS)*.
- Guo, Y., Codella, N. C., Karlinsky, L., Codella, J. V., Smith, J. R., Saenko, K., Rosing, T., & Feris, R. (2020). A broader study of cross-domain few-shot learning. *European Conf. Comput. Vis. (ECCV)*.
- Han, C., Murao, K., Noguchi, T., Kawata, Y., Uchiyama, F., Rundo, L., Nakayama, H., & Satoh, S. (2019). Learning more with less: Conditional PGGAN-based data augmentation for brain metastases detection using highly-rough annotation on MR images. *ACM Inter. Conf. Inf. Knowl. Manag. (CIKM)*, 119–127.
- Han, C., Rundo, L., Araki, R., Furukawa, Y., Mauri, G., Nakayama, H., & Hayashi, H. (2020). Infinite brain MR images: PGGAN-based data augmentation for tumor detection. *Neural approaches to dynamics of signal exchanges* (pp. 291–303). Springer.
- Hara, K., Kataoka, H., & Satoh, Y. (2018). Can spatiotemporal 3D CNNs retrace the history of 2D CNNs and ImageNet? *IEEE/CVF Inter. Conf. Comput. Vis. Pattern Recog. (CVPR)*.
- Hasani, B., & Mahoor, M. H. (2017a). Facial expression recognition using enhanced deep 3d convolutional neural networks. *IEEE Inter. Conf. Comput. Vis., Pattern Recog. Wksp. (CVPRW)*.
- Hasani, B., & Mahoor, M. H. (2017b). Spatio-temporal facial expression recognition using convolutional neural networks and conditional random fields. *IEEE Inter. Conf. Automat. Face Gesture Recog.*, 790–795.
- He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. *IEEE/CVF Inter. Conf. Comput. Vis. Pattern Recog. (CVPR)*.
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Comput.*, *9*(8), 1735–1780.
- Huang, G., Liu, Z., van der Maaten, L., & Weinberger, K. Q. (2017). Densely connected convolutional networks. *IEEE/CVF Inter. Conf. Comput. Vis. Pattern Recog. (CVPR)*.
- Huang, Y., Chen, F., Lv, S., & Wang, X. (2019). Facial expression recognition: A survey. *Sym.*, *11*(10).
- Huynh-The, T., & Kim, D.-S. (2019). Data augmentation for CNN-based 3D action recognition on small-scale datasets. *IEEE Inter. Conf. Ind. Inf. (INDIN)*.
- Jaderberg, M., Simonyan, K., Vedaldi, A., & Zisserman, A. (2016). Reading text in the wild with convolutional neural networks. *Inter. J. Comput. Vis.*, *116*(1), 1–20.
- Jena, R., Halder, S. S., & Sycara, K. (2020). MA3: Model agnostic adversarial augmentation for few shot learning. *IEEE Inter. Conf. Comput. Vis., Pattern Recog. Wksp. (CVPRW)*.
- Jung, H., Lee, S., Yim, J., Park, S., & Kim, J. (2015). Joint fine-tuning in deep neural networks for facial expression recognition. *IEEE Inter. Conf. Comput. Vis. (ICCV)*, 2983–2991.
- Kanade, T., Cohn, J. F., & Yingli Tian. (2000). Comprehensive database for facial expression analysis. *IEEE Inter. Conf. Automat. Face Gesture Recog.*, 46–53.
- Karras, T., Aittala, M., Hellsten, J., Laine, S., Lehtinen, J., & Aila, T. (2020). Training generative adversarial networks with limited data. *arXiv preprint arXiv:2006.06676*.
- Kaya, H., Gürpınar, F., & Salah, A. A. (2017). Video-based emotion recognition in the wild using deep transfer learning and score fusion [Multimodal Sentiment Analysis and Mining in the Wild Image and Vision Computing]. *Image Vis. Comput.*, *65*, 66–75.
- Kim, B.-K., Roh, J., Dong, S.-Y., & Lee, S.-Y. (2016). Hierarchical committee of deep convolutional neural networks for robust facial expression recognition. *J. Multi. User Inter.*, *10*(2), 173–189.
- Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

- Klaser, A., Marszałek, M., & Schmid, C. (2008). A spatio-temporal descriptor based on 3d-gradients. *British Mach. Vis. Conf. (BMVC)*, 275–1.
- Kuo, C.-M., Lai, S.-H., & Sarkis, M. (2018). A compact deep learning model for robust facial expression recognition. *IEEE Inter. Conf. Comput. Vis., Pattern Recog. Wksp. (CVPRW)*.
- Li, S., & Deng, W. (2018). Deep facial expression recognition: A survey. *IEEE Trans. Affect. Comput.*
- Liang, D., Liang, H., Yu, Z., & Zhang, Y. (2019). Deep convolutional BiLSTM fusion network for facial expression recognition. *Inter. J. Comput. Graphics*.
- Liu, M., Shan, S., Wang, R., & Chen, X. (2014). Learning expressionlets on spatio-temporal manifold for dynamic facial expression recognition. *IEEE/CVF Inter. Conf. Comput. Vis. Pattern Recog. (CVPR)*, 1749–1756.
- Liu, M., Li, S., Shan, S., Wang, R., & Chen, X. (2015). Deeply learning deformable facial action parts model for dynamic expression analysis. In D. Cremers, I. Reid, H. Saito, & M.-H. Yang (Eds.), *Asian conf. comput. vis. (ACCV)* (pp. 143–157). Springer International Publishing.
- Lopes, A. T., de Aguiar, E., Souza, A. F. D., & Oliveira-Santos, T. (2017). Facial expression recognition with convolutional neural networks: Coping with few data and the training sample order. *Pattern Recog.*, 61, 610–628.
- Lu, J., Gong, P., Ye, J., & Zhang, C. (2020). Learning from very few samples: A survey. *arXiv preprint arXiv:2009.02653*.
- Lucey, P., Cohn, J. F., Kanade, T., Saragih, J., Ambadar, Z., & Matthews, I. (2010). The extended cohn-kanade dataset (CK+): A complete dataset for action unit and emotion-specified expression. *IEEE/CVF Inter. Conf. Comput. Vis. Pattern Recog. (CVPR)*, 94–101.
- Maaten, L., Chen, M., Tyree, S., & Weinberger, K. (2013). Learning with marginalized corrupted features. *Inter. Conf. Mach. Learn. (ICML)*, 410–418.
- Marrero Fernandez, P. D., Guerrero Pena, F. A., Ren, T., & Cunha, A. (2019). FERAtt: Facial expression recognition with attention net. *IEEE Inter. Conf. Comput. Vis., Pattern Recog. Wksp. (CVPRW)*.
- Menon, L. T., Laurensi, I. A., Penna, M. C., Oliveira, L. E. S., & Britto, A. S. (2019). Data augmentation and transfer learning applied to charcoal image classification. *Inter. Conf. Sys. Signals Imag. Process. (IWSSIP)*, 69–74.
- Milicevic, M., Zubrinic, K., Obradovic, I., & Sjekavica, T. (2018). Data augmentation and transfer learning for limited dataset ship classification. *WSEAS Trans. Sys. Control*, 13, 460.
- Mollahosseini, A., Chan, D., & Mahoor, M. H. (2016). Going deeper in facial expression recognition using deep neural networks. *IEEE Wint. Conf. Appl. Comput. Vis. (WACV)*, 1–10.
- Nirkin, Y., Keller, Y., & Hassner, T. (2019). FSGAN: Subject agnostic face swapping and reenactment. *IEEE Inter. Conf. Comput. Vis. (ICCV)*.
- Pantic, M., Valstar, M., Rademaker, R., & Maat, L. (2005). Web-based database for facial expression analysis. *IEEE Inter. Conf. Multi. Expo. (ICME)*.
- Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., & Lerer, A. (2017). Automatic differentiation in PyTorch. *Wksp. Adv. Neural Inf. Process. Sys. (NeurIPS)*.
- Perez, L., & Wang, J. (2017). The effectiveness of data augmentation in image classification using deep learning. *arXiv preprint arXiv:1712.04621*, 11.
- Pinetz, T., Ruisz, J., & Soukup, D. (2019). Actual impact of gan augmentation on cnn classification performance. *Inter. Conf. Pattern Recog. Appl. Methods (ICPRAM)*, 15–23.
- Ramírez Rivera, A., & Chae, O. (2015). Spatiotemporal directional number transitional graph for dynamic texture recognition. *IEEE Trans. Pattern Anal. Mach. Intell.*, 37(10), 2146–2152.
- Ramírez Rivera, A., Rojas Castillo, J., & Chae, O. (2013). Local directional number pattern for face analysis: Face and expression recognition. *IEEE Trans. Image Process.*, 22(5), 1740–1752.
- Ramírez Rivera, A., Rojas Castillo, J., & Chae, O. (2015). Local directional texture pattern image descriptor. *Pattern Recog. Lett.*, 51(0), 94–100.
- Ratner, A. J., Ehrenberg, H., Hussain, Z., Dunnmon, J., & Ré, C. (2017). Learning to compose domain-specific transformations for data augmentation. *Adv. Neural Inf. Process. Sys. (NeurIPS)*, 3236–3246.
- Ryu, B., Ramírez Rivera, A., Kim, J., & Chae, O. (2017). Local directional ternary pattern for facial expression recognition. *IEEE Trans. Image Process.*, 26(12), 6006–6018.
- Shijie, J., Ping, W., Peiyi, J., & Siping, H. (2017). Research on data augmentation for image classification based on convolution neural networks. *Chinese Auto. Congr. (CAC)*, 4165–4170.
- Shin, H.-C., Roth, H. R., Gao, M., Lu, L., Xu, Z., Noguees, I., Yao, J., Mollura, D., & Summers, R. M. (2016). Deep convolutional neural networks for computer-aided detection: Cnn architectures, dataset characteristics and transfer learning. *IEEE Trans. Med. Imag.*, 35(5), 1285–1298.
- Shorten, C., & Khoshgoftaar, T. M. (2019). A survey on image data augmentation for deep learning. *J. Big Data*, 6(1), 60.
- Siarohin, A., Lathuilière, S., Tulyakov, S., Ricci, E., & Sebe, N. (2019). Animating arbitrary objects via deep motion transfer. *IEEE/CVF Inter. Conf. Comput. Vis. Pattern Recog. (CVPR)*.

- Simonyan, K., & Zisserman, A. (2015). Very deep convolutional networks for large-scale image recognition. *Inter. Conf. Learn. Represent. (ICLR)*.
- Soekhoe, D., van der Putten, P., & Plaat, A. (2016). On the impact of data set size in transfer learning using deep neural networks. *Inter. Symp. Intell. Data Anal. (IDA)*.
- Soomro, K., Zamir, A. R., & Shah, M. (2012). UCF101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*.
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: A simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.*, 15, 1929–1958.
- Srivastava, R. K., Greff, K., & Schmidhuber, J. (2015). Training very deep networks. *Adv. Neural Inf. Process. Sys. (NeurIPS)*, 2377–2385.
- Szegedy, C., Ioffe, S., & Vanhoucke, V. (2016). Inception-v4, inception-resnet and the impact of residual connections on learning. *AAAI Conf. Artif. Intell. (AAAI)*.
- Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., & Rabinovich, A. (2015). Going deeper with convolutions. *IEEE/CVF Inter. Conf. Comput. Vis. Pattern Recog. (CVPR)*, 1–9.
- Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., & Wojna, Z. (2016). Rethinking the inception architecture for computer vision. *IEEE/CVF Inter. Conf. Comput. Vis. Pattern Recog. (CVPR)*.
- Tajbakhsh, N., Shin, J. Y., Gurudu, S. R., Hurst, R. T., Kendall, C. B., Gotway, M. B., & Liang, J. (2016). Convolutional neural networks for medical image analysis: Full training or fine tuning? *IEEE Trans. Med. Imag.*, 35(5), 1299–1312.
- Tran, D., Bourdev, L., Fergus, R., Torresani, L., & Paluri, M. (2015). Learning spatiotemporal features with 3d convolutional networks. *IEEE Inter. Conf. Comput. Vis. (ICCV)*.
- Valstar, M., & Pantic, M. (2010). Induced disgust, happiness and surprise: An addition to the mmi facial expression database. *Inter. Wksp. Emotion*, 65.
- Wang, H., Gouk, H., Frank, E., Pfahringer, B., & Mayo, M. (2020). A comparison of machine learning methods for cross-domain few-shot learning. *Australasian Joint Conf. Artif. Intell. (AJCAI)*.
- Wang, X., Wang, K., & Lian, S. (2020). A survey on face data augmentation for the training of deep neural networks. *Neural Comput. Appls.*
- Wang, Y., Pan, X., Song, S., Zhang, H., Wu, C., & Huang, G. (2019). Implicit semantic data augmentation for deep networks. *Adv. Neural Inf. Process. Sys. (NeurIPS)*.
- Wu, W., Zhang, Y., Li, C., Qian, C., & Change Loy, C. (2018). ReenactGAN: Learning to reenact faces via boundary transfer. *European Conf. Comput. Vis. (ECCV)*.
- Yan, H. (2018). Collaborative discriminative multi-metric learning for facial expression recognition in video [Distance Metric Learning for Pattern Recognition]. *Pattern Recog.*, 75, 33–40.
- Zakharov, E., Shysheya, A., Burkov, E., & Lempitsky, V. (2019). Few-shot adversarial learning of realistic neural talking head models. *arXiv preprint arXiv:1905.08233*.
- Zhang, K., Huang, Y., Du, Y., & Wang, L. (2017). Facial expression recognition based on deep evolutionary spatial-temporal networks. *IEEE Trans. Image Process.*, 26(9), 4193–4203.
- Zhang, R., Che, T., Ghahramani, Z., Bengio, Y., & Song, Y. (2018). MetaGAN: An adversarial approach to few-shot learning. *Adv. Neural Inf. Process. Sys. (NeurIPS)*.
- Zhang, S., Pan, X., Cui, Y., Zhao, X., & Liu, L. (2019). Learning affective video features for facial expression recognition via hybrid deep learning. *IEEE Access*, 7, 32297–32304.
- Zhang, T., Zheng, W., Cui, Z., Zong, Y., & Li, Y. (2019). Spatial-temporal recurrent neural network for emotion recognition. *IEEE Trans. Cybern.*, 49(3), 839–847.
- Zhang, X., Wang, Z., Liu, D., & Ling, Q. (2019). DADA: Deep adversarial data augmentation for extremely low data regime classification. *IEEE Inter. Conf. Acoust., Speech, Signal Process. (ICASSP)*.
- Zhang, Y., Jia, G., Chen, L., Zhang, M., & Yong, J. (2019). Self-paced video data augmentation with dynamic images generated by generative adversarial networks. *arXiv preprint arXiv:1909.12929*.
- Zhao, G., Huang, X., Taini, M., Li, S. Z., & Pietikäinen, M. (2011). Facial expression recognition from near-infrared videos. *Image Vis. Comput.*, 29(9), 607–619.
- Zhao, G., & Pietikäinen, M. (2007). Dynamic texture recognition using local binary patterns with an application to facial expressions. *IEEE Trans. Pattern Anal. Mach. Intell.*, (6), 915–928.
- Zhao, L., Peng, X., Tian, Y., Kapadia, M., & Metaxas, D. (2018). Learning to forecast and refine residual motion for image-to-video generation. *European Conf. Comput. Vis. (ECCV)*.
- Zhao, M., Cong, Y., & Carin, L. (2020). On leveraging pretrained GANs for limited-data generation. *Inter. Conf. Mach. Learn. (ICML)*.
- Zhao, X., Liang, X., Liu, L., Li, T., Han, Y., Vasconcelos, N., & Yan, S. (2016). Peak-piloted deep network for facial expression recognition. In B. Leibe, J. Matas, N. Sebe, & M. Welling (Eds.), *European conf. comput. vis. (ECCV)* (pp. 425–442). Springer International Publishing.

Table A.1. Architecture details for VGG16 (Simonyan and Zisserman, 2015) plus LSTM.

Layer Name	Output Size	Configuration
Conv1	$224 \times 224 \times 64$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$
Max pool	$112 \times 112 \times 64$	2×2 max pool, stride 2
Conv2	$112 \times 112 \times 128$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$
Max pool	$56 \times 56 \times 128$	2×2 max pool, stride 2
Conv3	$56 \times 56 \times 256$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$
Max pool	$28 \times 28 \times 256$	2×2 max pool, stride 2
Conv4	$28 \times 28 \times 512$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$
Max pool	$14 \times 14 \times 512$	2×2 max pool, stride 2
Conv5	$14 \times 14 \times 512$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$
Max pool	$7 \times 7 \times 512$	2×2 max pool, stride 2
LSTM	25088×1024	1 Layer
Fully connected	512×256	512×256 fully connections
Fully connected	256×128	256×128 fully connections
Fully connected	$128 \times N_c$	$512 \times N_c$ fully connections
Softmax	N_c	N_c : number of classes

A. Models' Configuration

In the tables of this section, we show the implementation details of each model used in the experiments. The tables contain three columns: the name of the layer (or block), the output size, and the filter's configuration.

It should be noted that the configurations of models are the same as those introduced by the original authors; the only thing that changes is the elimination of the last layer (Softmax). Then, in 2D convolutional-based networks, an LSTM recurrent block is added to the last layer, which has a depth-one layer. Finally, in all models, a classifier block is stacked, composed of three dense layers that end in a softmax to make the prediction.

B. Single-Dataset k -Fold Cross-Validation Ablation Results

We show the detailed results of the experiments of all the methods when trained and tested on the same database in Table B.1. The configurations of each method correspond to it with and without fine-tuning and with different data augmentation.

We include results on the AFEW dataset (introduced in Section 4.3) for the single-, merged-, and cross-dataset scenarios. However, it is essential to remark that we kept the original partition of this dataset and did not perform 5-fold cross-validation experiments. The results reported in this and the following sections involving AFEW are not directly comparable to the other datasets for the single- and merged-dataset experiments.

C. Merged-Dataset k -Fold Cross-Validation Ablation Results

We show detailed results of the experiments of all the methods when trained on the merged folds of all databases and evaluated in each database individually in Table C.1. The configurations of each method correspond to it with and without fine-tuning and with different data augmentation.

Table A.2. Architecture details for Inception V3 (Szegedy, Ioffe, et al., 2016; Szegedy, W. Liu, et al., 2015; Szegedy, Vanhoucke, et al., 2016) plus LSTM.

Layer Name	Output Size	Configuration
Steam	$35 \times 35 \times 288$	$\begin{bmatrix} 3 \times 3, 32 \\ 3 \times 3, 32 \\ 3 \times 3, 64 \\ 3 \times 3, 80 \\ 3 \times 3, 192 \\ 3 \times 3, 288 \end{bmatrix} \times 1$
Inception A	$17 \times 17 \times 768$	$\begin{bmatrix} 1 \times 1 & 1 \times 1 & \text{Pool} & 1 \times 1 \\ 3 \times 3 & 3 \times 3 & 1 \times 1 \\ 3 \times 3 \\ & & \text{Concat} \end{bmatrix} \times 3$
Inception B	$8 \times 8 \times 1280$	$\begin{bmatrix} 1 \times 1 & 1 \times 1 & \text{Pool} & 1 \times 1 \\ 1 \times 7 & 1 \times 7 & 1 \times 1 \\ 7 \times 1 & 7 \times 1 \\ 1 \times 7 \\ 7 \times 1 \\ & & \text{Concat} \end{bmatrix} \times 3$
Inception C	$8 \times 8 \times 2048$	$\begin{bmatrix} & 1 \times 1 & & 1 \times 1 & \text{Pool} & 1 \times 1 \\ & 3 \times 3 & & 1 \times 3 & 3 \times 1 & 1 \times 1 \\ 1 \times 3 & 3 \times 1 & & & & \\ & & & & & \text{Concat} \end{bmatrix} \times 2$
Max Pool	$1 \times 1 \times 2048$	8×8 Max pool
LSTM	2048×512	1 Layer
Fully connected	512×256	73728×256 fully connections
Fully connected	256×128	256×128 fully connections
Fully connected	$128 \times N_c$	$512 \times N_c$ fully connections
Softmax	N_c	N_c : number of classes

D. Cross-Dataset Versus All Results

Each table shows the results of a model-initialization-augmentation combination. Due to time and computational resource constraints, we decided to experiment with fine-tuning as our initialization scheme since the random one yielded lower performance in comparison.

E. Statistical Hypothesis Tests for Experiments

In order to compare the results of our different experiments, we performed a Wilcoxon Signed-Rank Test. As suggested in the literature (Demšar, 2006; Dietterich, 1998), more thorough comparisons are needed to conclude the similarity between methods. However, this comparison is commonly ignored due to the lack of information on used folds in different databases.

Since we are constructing our experiments, we took care of using the same folds for all experiments. Hence, we can consider two methods running on the same folds a paired experiment, and, thus, we used the Wilcoxon Signed-Rank test (a non-parametric test for comparison) (Demšar, 2006).

In the the tables of this section we present the p -value of each test pair-wise per database for the single- and merged-database (as explained in Sections 4.1 and 4.2, respectively). We show the statistical test at two- and one-tail. The former evaluates whether the methods' accuracy is equal, while the former evaluate if one is greater. We denote by shading the p -values the rejection of the null hypothesis for each test.

Table A.3. Architecture details for ResNet-18 (He et al., 2016) plus LSTM.

Layer Name	Output Size	Configuration
Conv1	$112 \times 112 \times 64$	$7 \times 7, 64$, stride 2
Max pool	$56 \times 56 \times 64$	$2 \times$ max pool, stride 2
Conv2	$56 \times 56 \times 64$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$
Conv3	$28 \times 28 \times 128$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$
Conv4	$14 \times 14 \times 256$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$
Conv5	$7 \times 7 \times 512$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$
Average pool	$1 \times 1 \times 512$	7×7 average pool
LSTM	512×512	1 Layer
Fully connected	512×256	512×256 fully connections
Fully connected	256×128	256×128 fully connections
Fully connected	$128 \times N_c$	$512 \times N_c$ fully connections
Softmax	N_c	N_c : number of classes

Table A.4. Architecture details for ResNet-101 (He et al., 2016) plus LSTM.

Layer Name	Output Size	Configuration
Conv1	$112 \times 112 \times 64$	$7 \times 7, 64$, stride 2
Max pool	$56 \times 56 \times 64$	2×2 max pool, stride 2
Conv2	$56 \times 56 \times 256$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
Conv3	$28 \times 28 \times 512$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$
Conv4	$14 \times 14 \times 1024$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$
Conv5	$7 \times 7 \times 2048$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
Average pool	$1 \times 1 \times 2048$	7×7 average pool
LSTM	2048×512	1 Layer
Fully connected	512×256	512×256 fully connections
Fully connected	256×128	256×128 fully connections
Fully connected	$128 \times N_c$	$512 \times N_c$ fully connections
Softmax	N_c	N_c : number of classes

Table A.5. Architecture details for C3D (Tran et al., 2015).

Layer Name	Output Size	Configuration
Conv1	$100 \times 100 \times 16 \times 64$	$[3 \times 3 \times 3, 64] \times 1$
Max pool	$50 \times 50 \times 16 \times 64$	$2 \times 2 \times 1$ max pool, stride 2
Conv2	$50 \times 50 \times 16 \times 128$	$[3 \times 3 \times 3, 128] \times 1$
Max pool	$25 \times 25 \times 16 \times 128$	$2 \times 2 \times 1$ max pool, stride 2
Conv3	$25 \times 25 \times 16 \times 256$	$[3 \times 3 \times 3, 256] \times 2$
Max pool	$12 \times 12 \times 16 \times 256$	$2 \times 2 \times 1$ max pool, stride 2
Conv4	$12 \times 12 \times 16 \times 512$	$[3 \times 3 \times 3, 512] \times 2$
Max pool	$6 \times 6 \times 16 \times 512$	$2 \times 2 \times 1$ max pool, stride 2
Conv5	$6 \times 6 \times 16 \times 512$	$[3 \times 3 \times 3, 512] \times 2$
Max pool	$3 \times 3 \times 16 \times 512$	$2 \times 2 \times 1$ max pool, stride 2
Fully connected	73728×256	73728×256 fully connections
Fully connected	256×128	256×128 fully connections
Fully connected	$128 \times N_c$	$512 \times N_c$ fully connections
Softmax	N_c	N_c : number of classes

Table A.6. Architecture details for I3D (Carreira and Zisserman, 2017).

Layer Name	Output Size	Configuration
Conv 1	$112 \times 112 \times 32 \times 64$	$[7 \times 7 \times 7, 64] \times 1$
Max Pool	$56 \times 56 \times 16 \times 64$	$2 \times 2 \times 2$ Max pool, stride 2
Conv 2	$56 \times 56 \times 16 \times 192$	$[3 \times 3 \times 3, 192] \times 1$
Max Pool	$28 \times 28 \times 8 \times 192$	$2 \times 2 \times 2$ Max pool, stride 2
Inception 3a	$28 \times 28 \times 8 \times 256$	$\begin{bmatrix} 1 \times 1 \times 1 & 1 \times 1 \times 1 & 1 \times 1 \times 1 & \text{Max Pool} \\ & 3 \times 3 \times 3 & 3 \times 3 \times 3 & 1 \times 1 \times 1 \\ & & \text{Concat} & \end{bmatrix} \times 1$
Inception 3b	$28 \times 28 \times 8 \times 480$	$\begin{bmatrix} 1 \times 1 \times 1 & 1 \times 1 \times 1 & 1 \times 1 \times 1 & \text{Max Pool} \\ & 3 \times 3 \times 3 & 3 \times 3 \times 3 & 1 \times 1 \times 1 \\ & & \text{Concat} & \end{bmatrix} \times 1$
Max Pool	$14 \times 14 \times 4 \times 480$	$2 \times 2 \times 2$ Max pool, stride 2
Inception 4a	$14 \times 14 \times 4 \times 512$	$\begin{bmatrix} 1 \times 1 \times 1 & 1 \times 1 \times 1 & 1 \times 1 \times 1 & \text{Max Pool} \\ & 3 \times 3 \times 3 & 3 \times 3 \times 3 & 1 \times 1 \times 1 \\ & & \text{Concat} & \end{bmatrix} \times 1$
Inception 4b	$14 \times 14 \times 4 \times 512$	$\begin{bmatrix} 1 \times 1 \times 1 & 1 \times 1 \times 1 & 1 \times 1 \times 1 & \text{Max Pool} \\ & 3 \times 3 \times 3 & 3 \times 3 \times 3 & 1 \times 1 \times 1 \\ & & \text{Concat} & \end{bmatrix} \times 1$
Inception 4c	$14 \times 14 \times 4 \times 512$	$\begin{bmatrix} 1 \times 1 \times 1 & 1 \times 1 \times 1 & 1 \times 1 \times 1 & \text{Max Pool} \\ & 3 \times 3 \times 3 & 3 \times 3 \times 3 & 1 \times 1 \times 1 \\ & & \text{Concat} & \end{bmatrix} \times 1$
Inception 4d	$14 \times 14 \times 4 \times 528$	$\begin{bmatrix} 1 \times 1 \times 1 & 1 \times 1 \times 1 & 1 \times 1 \times 1 & \text{Max Pool} \\ & 3 \times 3 \times 3 & 3 \times 3 \times 3 & 1 \times 1 \times 1 \\ & & \text{Concat} & \end{bmatrix} \times 1$
Inception 4e	$14 \times 14 \times 4 \times 832$	$\begin{bmatrix} 1 \times 1 \times 1 & 1 \times 1 \times 1 & 1 \times 1 \times 1 & \text{Max Pool} \\ & 3 \times 3 \times 3 & 3 \times 3 \times 3 & 1 \times 1 \times 1 \\ & & \text{Concat} & \end{bmatrix} \times 1$
Max Pool	$7 \times 7 \times 2 \times 832$	$2 \times 2 \times 2$ Max pool, stride 2
Inception 5a	$7 \times 7 \times 2 \times 832$	$\begin{bmatrix} 1 \times 1 \times 1 & 1 \times 1 \times 1 & 1 \times 1 \times 1 & \text{Max Pool} \\ & 3 \times 3 \times 3 & 3 \times 3 \times 3 & 1 \times 1 \times 1 \\ & & \text{Concat} & \end{bmatrix} \times 1$
Inception 5b	$7 \times 7 \times 2 \times 1024$	$\begin{bmatrix} 1 \times 1 \times 1 & 1 \times 1 \times 1 & 1 \times 1 \times 1 & \text{Max Pool} \\ & 3 \times 3 \times 3 & 3 \times 3 \times 3 & 1 \times 1 \times 1 \\ & & \text{Concat} & \end{bmatrix} \times 1$
Average Pool	$1 \times 1 \times 1 \times 1024$	$7 \times 7 \times 2$ Max pool
Fully connected	1024×256	73728×256 fully connections
Fully connected	256×128	256×128 fully connections
Fully connected	$128 \times N_c$	$512 \times N_c$ fully connections
Softmax	N_c	N_c : number of classes

Table A.7. Architecture details for ResNet3D-18 (Hara et al., 2018).

Layer Name	Output Size	Configuration
Conv1	$100 \times 100 \times 25 \times 64$	$5 \times 5 \times 5, 64$, stride 2
Max3D pool	$50 \times 50 \times 12 \times 64$	$2 \times 2 \times 2$ max pool, stride 2
Conv2	$50 \times 50 \times 12 \times 64$	$\begin{bmatrix} 3 \times 3 \times 1, 64 \\ 3 \times 3 \times 1, 64 \end{bmatrix} \times 2$
Conv3	$25 \times 25 \times 6 \times 128$	$\begin{bmatrix} 3 \times 3 \times 1, 128 \\ 3 \times 3 \times 1, 128 \end{bmatrix} \times 2$
Conv4	$12 \times 12 \times 3 \times 256$	$\begin{bmatrix} 3 \times 3 \times 3, 256 \\ 3 \times 3 \times 3, 256 \end{bmatrix} \times 2$
Conv5	$6 \times 6 \times 1 \times 512$	$\begin{bmatrix} 3 \times 3 \times 3, 512 \\ 3 \times 3 \times 3, 512 \end{bmatrix} \times 2$
Average3D pool	$1 \times 1 \times 1 \times 512$	6×6 average pool
Fully connected	512×256	512×256 fully connections
Fully connected	256×128	256×128 fully connections
Fully connected	$128 \times N_c$	$512 \times N_c$ fully connections
Softmax	N_c	N_c : number of classes

Table A.8. Architecture details for ResNet3D-101 (Hara et al., 2018).

Layer Name	Output Size	Configuration
Conv1	$100 \times 100 \times 25 \times 64$	$[7 \times 7 \times 7, 64] \times 1$
Max3D pool	$50 \times 50 \times 12 \times 64$	$2 \times 2 \times 2$ max pool, stride 2
Conv2	$50 \times 50 \times 6 \times 256$	$\begin{bmatrix} 1 \times 1 \times 1, 64 \\ 3 \times 3 \times 3, 64 \\ 1 \times 1 \times 1, 256 \end{bmatrix} \times 3$
Conv3	$25 \times 25 \times 3 \times 512$	$\begin{bmatrix} 1 \times 1 \times 1, 128 \\ 3 \times 3 \times 3, 128 \\ 1 \times 1 \times 1, 512 \end{bmatrix} \times 4$
Conv4	$12 \times 12 \times 1 \times 1024$	$\begin{bmatrix} 1 \times 1 \times 1, 256 \\ 3 \times 3 \times 3, 256 \\ 1 \times 1 \times 1, 1024 \end{bmatrix} \times 23$
Conv5	$6 \times 6 \times 1 \times 2048$	$\begin{bmatrix} 1 \times 1 \times 1, 512 \\ 3 \times 3 \times 3, 512 \\ 1 \times 1 \times 1, 2048 \end{bmatrix} \times 3$
Average3D pool	$1 \times 1 \times 1 \times 2048$	$6 \times 6 \times 1$ average pool
Fully connected	2048×256	512×256 fully connections
Fully connected	256×128	256×128 fully connections
Fully connected	$128 \times N_c$	$512 \times N_c$ fully connections
Softmax	N_c	N_c : number of classes

Table A.9. Architecture details for C3D-block-LSTM.

Layer Name	Output Size	Configuration
Conv1	$100 \times 100 \times 16 \times 64$	$[3 \times 3 \times 3, 64] \times 1$
Max pool	$50 \times 50 \times 16 \times 64$	$2 \times 2 \times 1$ max pool, stride 2
Conv2	$50 \times 50 \times 16 \times 128$	$[3 \times 3 \times 3, 128] \times 1$
Max pool	$25 \times 25 \times 16 \times 128$	$2 \times 2 \times 1$ max pool, stride 2
Conv3	$25 \times 25 \times 16 \times 256$	$[3 \times 3 \times 3, 256] \times 2$
Max pool	$12 \times 12 \times 16 \times 256$	$2 \times 2 \times 1$ max pool, stride 2
Conv4	$12 \times 12 \times 16 \times 512$	$[3 \times 3 \times 3, 512] \times 2$
Max pool	$6 \times 6 \times 16 \times 512$	$2 \times 2 \times 1$ max pool, stride 2
Conv5	$6 \times 6 \times 16 \times 512$	$[3 \times 3 \times 3, 512] \times 2$
Max pool	$3 \times 3 \times 16 \times 512$	$2 \times 2 \times 1$ max pool, stride 2
LSTM	73728×512	1 Layer
Fully connected	512×256	73728×256 fully connections
Fully connected	256×128	256×128 fully connections
Fully connected	$128 \times N_c$	$512 \times N_c$ fully connections
Softmax	N_c	N_c : number of classes

Table B.1. Ablation study of several methods when trained and tested on the **same database** with random initialization (RI) or fine-tuning (FT), i.e., pre-training and then transfer learning, and with data augmentation using random transformations (DA), using synthetic generated data (SD), and without it (-).

Architecture	Techniques			Datasets			
	Init.	Data Aug.	CK+	MMI	OULU	MUG	AFEW
VGG16-LSTM	RI	-	0.847 ± 0.049	0.202 ± 0.079	0.656 ± 0.037	0.186 ± 0.006	0.188
	RI	DA	0.655 ± 0.092	0.194 ± 0.058	0.450 ± 0.077	0.175 ± 0.048	0.199
	RI	SD	0.141 ± 0.085	0.155 ± 0.024	0.128 ± 0.031	0.568 ± 0.334	0.165
	FT	-	0.911 ± 0.040	0.612 ± 0.044	0.431 ± 0.324	0.862 ± 0.041	0.270
	FT	DA	0.954 ± 0.032	0.582 ± 0.045	0.750 ± 0.039	0.882 ± 0.035	0.225
	FT	SD	0.920 ± 0.033	0.547 ± 0.047	0.752 ± 0.066	0.840 ± 0.051	0.262
	InceptionV3-LSTM	RI	-	0.732 ± 0.001	0.241 ± 0.002	0.163 ± 0.005	0.605 ± 0.145
RI		DA	0.321 ± 0.052	0.232 ± 0.042	0.184 ± 0.023	0.410 ± 0.172	0.162
RI		SD	0.622 ± 0.176	0.321 ± 0.036	0.448 ± 0.070	0.647 ± 0.166	0.162
FT		-	0.820 ± 0.056	0.475 ± 0.077	0.588 ± 0.035	0.874 ± 0.028	0.217
FT		DA	0.924 ± 0.019	0.353 ± 0.112	0.621 ± 0.073	0.809 ± 0.057	0.152
FT		SD	0.792 ± 0.090	0.425 ± 0.110	0.769 ± 0.045	0.813 ± 0.061	0.264
ResNet18-LSTM		RI	-	0.560 ± 0.089	0.224 ± 0.083	0.444 ± 0.035	0.687 ± 0.062
	RI	DA	0.627 ± 0.107	0.363 ± 0.102	0.508 ± 0.067	0.673 ± 0.051	0.154
	RI	SD	0.670 ± 0.113	0.335 ± 0.042	0.533 ± 0.044	0.771 ± 0.039	0.191
	FT	-	0.783 ± 0.090	0.479 ± 0.136	0.692 ± 0.042	0.867 ± 0.025	0.241
	FT	DA	0.771 ± 0.053	0.473 ± 0.120	0.648 ± 0.071	0.823 ± 0.041	0.220
	FT	SD	0.887 ± 0.054	0.412 ± 0.073	0.733 ± 0.040	0.866 ± 0.041	0.209
	ResNet101-LSTM	RI	-	0.644 ± 0.038	0.250 ± 0.085	0.465 ± 0.073	0.623 ± 0.063
RI		DA	0.542 ± 0.106	0.301 ± 0.107	0.419 ± 0.031	0.721 ± 0.001	0.002
RI		SD	0.654 ± 0.068	0.210 ± 0.070	0.596 ± 0.018	0.644 ± 0.070	0.173
FT		-	0.805 ± 0.013	0.480 ± 0.077	0.754 ± 0.036	0.851 ± 0.060	0.173
FT		DA	0.796 ± 0.130	0.463 ± 0.161	0.679 ± 0.046	0.804 ± 0.026	0.186
FT		SD	0.899 ± 0.029	0.326 ± 0.168	0.744 ± 0.067	0.740 ± 0.139	0.181
C3D		RI	-	0.254 ± 0.017	0.235 ± 0.050	0.172 ± 0.004	0.187 ± 0.005
	RI	DA	0.768 ± 0.026	0.280 ± 0.033	0.241 ± 0.029	0.190 ± 0.005	0.236
	RI	SD	0.154 ± 0.053	0.158 ± 0.082	0.139 ± 0.044	0.155 ± 0.030	0.165
	FT	-	0.878 ± 0.021	0.476 ± 0.073	0.648 ± 0.042	0.182 ± 0.011	0.165
	FT	DA	0.869 ± 0.065	0.398 ± 0.151	0.165 ± 0.004	0.188 ± 0.006	0.165
	FT	SD	0.859 ± 0.048	0.264 ± 0.182	0.250 ± 0.024	0.399 ± 0.292	0.230
	C3D-Block-LSTM	RI	-	0.713 ± 0.035	0.207 ± 0.075	0.175 ± 0.025	0.172 ± 0.032
RI		DA	0.254 ± 0.017	0.194 ± 0.082	0.205 ± 0.036	0.193 ± 0.041	0.165
RI		SD	0.149 ± 0.073	0.118 ± 0.040	0.135 ± 0.065	0.176 ± 0.022	0.165
FT		-	0.853 ± 0.027	0.486 ± 0.113	0.650 ± 0.030	0.861 ± 0.047	0.183
FT		DA	0.872 ± 0.069	0.549 ± 0.064	0.717 ± 0.034	0.854 ± 0.022	0.233
FT		SD	0.890 ± 0.038	0.491 ± 0.098	0.542 ± 0.194	0.678 ± 0.276	0.249
I3D		RI	-	0.767 ± 0.066	0.314 ± 0.086	0.796 ± 0.037	0.851 ± 0.018
	RI	DA	0.850 ± 0.046	0.409 ± 0.109	0.685 ± 0.056	0.851 ± 0.020	0.202
	RI	SD	0.808 ± 0.092	0.346 ± 0.075	0.671 ± 0.078	0.837 ± 0.028	0.221
	FT	-	0.912 ± 0.044	0.509 ± 0.083	0.771 ± 0.040	0.869 ± 0.020	0.283
	FT	DA	0.896 ± 0.050	0.530 ± 0.080	0.752 ± 0.045	0.866 ± 0.034	0.353
	FT	SD	0.878 ± 0.032	0.445 ± 0.086	0.710 ± 0.048	0.850 ± 0.043	0.401
	ResNet3D-18	RI	-	0.756 ± 0.045	0.324 ± 0.078	0.577 ± 0.022	0.848 ± 0.021
RI		DA	0.768 ± 0.051	0.312 ± 0.109	0.608 ± 0.036	0.782 ± 0.081	0.162
RI		SD	0.856 ± 0.040	0.319 ± 0.047	0.588 ± 0.068	0.802 ± 0.001	0.213
FT		-	0.804 ± 0.019	0.312 ± 0.129	0.646 ± 0.029	0.873 ± 0.036	0.188
FT		DA	0.832 ± 0.034	0.492 ± 0.086	0.748 ± 0.052	0.816 ± 0.051	0.207
FT		SD	0.859 ± 0.035	0.354 ± 0.089	0.696 ± 0.044	0.843 ± 0.040	0.233
ResNet3D-101		RI	-	0.713 ± 0.063	0.269 ± 0.071	0.537 ± 0.045	0.737 ± 0.049
	RI	DA	0.658 ± 0.075	0.242 ± 0.092	0.348 ± 0.008	0.803 ± 0.001	0.173
	RI	SD	0.774 ± 0.068	0.315 ± 0.052	0.481 ± 0.043	0.742 ± 0.072	0.188
	FT	-	0.920 ± 0.034	0.538 ± 0.110	0.781 ± 0.025	0.896 ± 0.019	0.220
	FT	DA	0.911 ± 0.042	0.585 ± 0.092	0.756 ± 0.027	0.850 ± 0.057	0.304
	FT	SD	0.896 ± 0.034	0.462 ± 0.100	0.756 ± 0.044	0.865 ± 0.034	0.330

Table C.1. Ablation study of several methods when trained and tested on the **combination of databases** with random initialization (RI) or fine-tuning (FT), i.e., pre-training and then transfer learning, and with data augmentation using random transformations (DA), using synthetic generated data (SD), and without it (-).

Architecture	Techniques			Datasets			
	Init.	Data Aug.	CK+	MMI	OULU	MUG	All
VGG16-LSTM	RI	-	0.282 ± 0.030	0.236 ± 0.045	0.183 ± 0.053	0.201 ± 0.001	0.226 ± 0.032
	RI	DA	0.249 ± 0.043	0.231 ± 0.035	0.193 ± 0.038	0.184 ± 0.012	0.214 ± 0.032
	RI	SD	0.267 ± 0.011	0.202 ± 0.078	0.207 ± 0.014	0.217 ± 0.007	0.223 ± 0.028
	FT	-	0.929 ± 0.014	0.676 ± 0.101	0.787 ± 0.019	0.903 ± 0.023	0.824 ± 0.039
	FT	DA	0.929 ± 0.045	0.419 ± 0.084	0.506 ± 0.058	0.596 ± 0.051	0.612 ± 0.060
	FT	SD	0.938 ± 0.034	0.688 ± 0.067	0.779 ± 0.033	0.909 ± 0.021	0.829 ± 0.039
InceptionV3-LSTM	RI	-	0.505 ± 0.117	0.327 ± 0.093	0.392 ± 0.039	0.565 ± 0.063	0.447 ± 0.078
	RI	DA	0.290 ± 0.091	0.207 ± 0.066	0.167 ± 0.007	0.191 ± 0.008	0.213 ± 0.043
	RI	SD	0.372 ± 0.143	0.244 ± 0.056	0.237 ± 0.056	0.507 ± 0.124	0.340 ± 0.095
	FT	-	0.886 ± 0.038	0.596 ± 0.211	0.654 ± 0.204	0.740 ± 0.256	0.719 ± 0.177
	FT	DA	0.549 ± 0.342	0.155 ± 0.055	0.181 ± 0.021	0.165 ± 0.027	0.263 ± 0.111
	FT	SD	0.854 ± 0.072	0.587 ± 0.129	0.704 ± 0.080	0.825 ± 0.076	0.742 ± 0.089
ResNet18-LSTM	RI	-	0.584 ± 0.069	0.407 ± 0.099	0.475 ± 0.048	0.620 ± 0.093	0.522 ± 0.077
	RI	DA	0.550 ± 0.245	0.196 ± 0.063	0.196 ± 0.040	0.198 ± 0.058	0.285 ± 0.102
	RI	SD	0.606 ± 0.116	0.387 ± 0.088	0.494 ± 0.053	0.599 ± 0.098	0.521 ± 0.089
	FT	-	0.856 ± 0.149	0.615 ± 0.066	0.683 ± 0.085	0.816 ± 0.101	0.742 ± 0.100
	FT	DA	0.776 ± 0.324	0.107 ± 0.036	0.181 ± 0.031	0.160 ± 0.021	0.306 ± 0.103
	FT	SD	0.866 ± 0.087	0.606 ± 0.087	0.717 ± 0.073	0.822 ± 0.059	0.753 ± 0.076
ResNet101-LSTM	RI	-	0.319 ± 0.155	0.186 ± 0.068	0.379 ± 0.124	0.523 ± 0.130	0.352 ± 0.119
	RI	DA	0.249 ± 0.134	0.187 ± 0.096	0.183 ± 0.028	0.172 ± 0.040	0.198 ± 0.075
	RI	SD	0.348 ± 0.173	0.245 ± 0.046	0.402 ± 0.148	0.535 ± 0.136	0.382 ± 0.126
	FT	-	0.783 ± 0.183	0.491 ± 0.143	0.525 ± 0.267	0.816 ± 0.054	0.654 ± 0.162
	FT	DA	0.351 ± 0.247	0.230 ± 0.049	0.165 ± 0.004	0.185 ± 0.007	0.233 ± 0.077
	FT	SD	0.761 ± 0.215	0.499 ± 0.150	0.579 ± 0.179	0.651 ± 0.258	0.623 ± 0.200
C3D	RI	-	0.254 ± 0.028	0.194 ± 0.049	0.192 ± 0.017	0.200 ± 0.021	0.210 ± 0.028
	RI	DA	0.253 ± 0.019	0.204 ± 0.058	0.184 ± 0.029	0.187 ± 0.014	0.207 ± 0.030
	RI	SD	0.268 ± 0.011	0.218 ± 0.034	0.207 ± 0.019	0.229 ± 0.020	0.231 ± 0.021
	FT	-	0.864 ± 0.022	0.570 ± 0.053	0.667 ± 0.055	0.835 ± 0.043	0.734 ± 0.043
	FT	DA	0.857 ± 0.036	0.324 ± 0.069	0.352 ± 0.041	0.446 ± 0.116	0.495 ± 0.066
	FT	SD	0.873 ± 0.031	0.583 ± 0.050	0.700 ± 0.067	0.843 ± 0.049	0.750 ± 0.049
C3D-Block-LSTM	RI	-	0.301 ± 0.003	0.204 ± 0.009	0.185 ± 0.031	0.161 ± 0.035	0.213 ± 0.019
	RI	DA	0.288 ± 0.003	0.237 ± 0.040	0.190 ± 0.004	0.165 ± 0.021	0.220 ± 0.017
	RI	SD	0.268 ± 0.011	0.207 ± 0.075	0.173 ± 0.013	0.183 ± 0.010	0.208 ± 0.027
	FT	-	0.890 ± 0.045	0.612 ± 0.051	0.690 ± 0.026	0.853 ± 0.052	0.761 ± 0.043
	FT	DA	0.769 ± 0.243	0.411 ± 0.124	0.365 ± 0.107	0.398 ± 0.112	0.486 ± 0.146
	FT	SD	0.886 ± 0.048	0.601 ± 0.045	0.713 ± 0.036	0.865 ± 0.040	0.766 ± 0.042
I3D	RI	-	0.821 ± 0.061	0.447 ± 0.040	0.708 ± 0.038	0.807 ± 0.024	0.696 ± 0.041
	RI	DA	0.667 ± 0.141	0.252 ± 0.060	0.244 ± 0.033	0.296 ± 0.071	0.365 ± 0.076
	RI	SD	0.759 ± 0.083	0.559 ± 0.091	0.640 ± 0.100	0.742 ± 0.083	0.675 ± 0.089
	FT	-	0.867 ± 0.029	0.635 ± 0.069	0.748 ± 0.039	0.853 ± 0.028	0.776 ± 0.041
	FT	DA	0.867 ± 0.013	0.196 ± 0.067	0.194 ± 0.045	0.172 ± 0.032	0.357 ± 0.039
	FT	SD	0.864 ± 0.034	0.642 ± 0.052	0.748 ± 0.039	0.854 ± 0.026	0.777 ± 0.038
ResNet3D-18	RI	-	0.822 ± 0.068	0.425 ± 0.070	0.602 ± 0.047	0.801 ± 0.049	0.663 ± 0.058
	RI	DA	0.807 ± 0.106	0.286 ± 0.048	0.310 ± 0.080	0.318 ± 0.065	0.430 ± 0.075
	RI	SD	0.759 ± 0.112	0.412 ± 0.054	0.546 ± 0.105	0.745 ± 0.078	0.615 ± 0.087
	FT	-	0.813 ± 0.041	0.479 ± 0.066	0.604 ± 0.043	0.793 ± 0.031	0.672 ± 0.045
	FT	DA	0.819 ± 0.044	0.194 ± 0.067	0.210 ± 0.043	0.196 ± 0.034	0.355 ± 0.047
	FT	SD	0.858 ± 0.027	0.517 ± 0.072	0.652 ± 0.044	0.822 ± 0.008	0.712 ± 0.038
ResNet3D-101	RI	-	0.486 ± 0.046	0.268 ± 0.087	0.354 ± 0.037	0.508 ± 0.041	0.404 ± 0.053
	RI	DA	0.453 ± 0.038	0.243 ± 0.056	0.215 ± 0.018	0.219 ± 0.045	0.282 ± 0.039
	RI	SD	0.515 ± 0.065	0.309 ± 0.071	0.356 ± 0.017	0.521 ± 0.092	0.425 ± 0.061
	FT	-	0.942 ± 0.014	0.646 ± 0.112	0.810 ± 0.062	0.897 ± 0.028	0.824 ± 0.054
	FT	DA	0.929 ± 0.039	0.362 ± 0.104	0.481 ± 0.115	0.517 ± 0.138	0.572 ± 0.099
	FT	SD	0.951 ± 0.015	0.659 ± 0.082	0.790 ± 0.038	0.884 ± 0.028	0.821 ± 0.041

Table D.1. Classification accuracies of VGG16-LSTM with Fine Tuning when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.405	0.563	0.541	0.169
MMI	0.786	–	0.531	0.688	0.192
OULU	0.738	0.321	–	0.406	0.136
MUG	0.790	0.524	0.465	–	0.177
AFEW	0.288	0.268	0.238	0.263	–

Table D.2. Classification accuracies of VGG16-LSTM with Fine Tuning + Data augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.405	0.565	0.574	0.173
MMI	0.650	–	0.410	0.622	0.168
OULU	0.812	0.506	–	0.522	0.175
MUG	0.754	0.524	0.492	–	0.207
AFEW	0.301	0.262	0.254	0.295	–

Table D.3. Classification accuracies of VGG16-LSTM with Fine Tuning + Synthetic augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.452	0.581	0.650	0.186
MMI	0.269	–	0.167	0.190	0.177
OULU	0.191	0.137	–	0.164	0.113
MUG	0.231	0.321	0.213	–	0.123
AFEW	0.146	0.167	0.167	0.180	–

Table D.4. Classification accuracies of VGG16-LSTM with Fine Tuning + Data augmentation + Synthetic augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.417	0.571	0.576	0.156
MMI	0.269	–	0.167	0.190	0.177
OULU	0.191	0.137	–	0.164	0.113
MUG	0.146	0.167	0.167	–	0.194
AFEW	0.146	0.167	0.167	0.180	–

Table D.5. Classification accuracies of InceptionV3-LSTM with Fine Tuning when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.167	0.167	0.217	0.120
MMI	0.139	–	0.173	0.371	0.188
OULU	0.699	0.345	–	0.238	0.182
MUG	0.146	0.351	0.171	–	0.204
AFEW	0.269	0.256	0.177	0.253	–

Table D.6. Classification accuracies of InceptionV3-LSTM with Fine Tuning + Data augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.196	0.213	0.244	0.140
MMI	0.485	–	0.269	0.575	0.135
OULU	0.718	0.363	–	0.377	0.170
MUG	0.521	0.405	0.338	–	0.169
AFEW	0.269	0.256	0.177	0.253	–

Table D.7. Classification accuracies of InceptionV3-LSTM with Fine Tuning + Synthetic augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.202	0.219	0.249	0.154
MMI	0.136	–	0.177	0.147	0.135
OULU	0.181	0.173	–	0.180	0.143
MUG	0.275	0.214	0.163	–	0.176
AFEW	0.091	0.167	0.167	0.143	–

Table D.8. Classification accuracies of InceptionV3-LSTM with Fine Tuning + Data augmentation + Synthetic augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.167	0.146	0.200	0.148
MMI	0.091	–	0.179	0.150	0.155
OULU	0.214	0.202	–	0.236	0.143
MUG	0.269	0.208	0.169	–	0.177
AFEW	0.091	0.173	0.179	0.143	–

Table D.9. Classification accuracies of ResNet18-LSTM with Fine Tuning when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.179	0.160	0.177	0.175
MMI	0.560	–	0.283	0.524	0.182
OULU	0.660	0.363	–	0.455	0.138
MUG	0.405	0.429	0.325	–	0.135
AFEW	0.269	0.256	0.223	0.276	–

Table D.10. Classification accuracies of ResNet18-LSTM with Fine Tuning + Data augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.244	0.183	0.244	0.176
MMI	0.599	–	0.381	0.651	0.149
OULU	0.660	0.345	–	0.488	0.151
MUG	0.780	0.619	0.473	–	0.148
AFEW	0.265	0.149	0.121	0.223	–

Table D.11. Classification accuracies of ResNet18-LSTM with Fine Tuning + Synthetic augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.137	0.169	0.164	0.110
MMI	0.081	–	0.167	0.137	0.124
OULU	0.136	0.167	–	0.182	0.171
MUG	0.269	0.125	0.150	–	0.156
AFEW	0.087	0.143	0.167	0.139	–

Table D.12. Classification accuracies of ResNet18-LSTM with Fine Tuning + Data augmentation + Synthetic augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.119	0.215	0.160	0.166
MMI	0.084	–	0.154	0.137	0.134
OULU	0.191	0.137	–	0.164	0.113
MUG	0.149	0.167	0.167	–	0.191
AFEW	0.081	0.143	0.167	0.137	–

Table D.13. Classification accuracies of ResNet101-LSTM with Fine Tuning when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.167	0.169	0.180	0.194
MMI	0.430	–	0.185	0.529	0.186
OULU	0.812	0.435	–	0.632	0.156
MUG	0.104	0.244	0.183	–	0.129
AFEW	0.340	0.238	0.198	0.301	–

Table D.14. Classification accuracies of ResNet101-LSTM with Fine Tuning + Data augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.208	0.171	0.190	0.129
MMI	0.657	–	0.402	0.654	0.156
OULU	0.874	0.494	–	0.606	0.165
MUG	0.447	0.393	0.325	–	0.161
AFEW	0.382	0.244	0.175	0.207	–

Table D.15. Classification accuracies of ResNet101-LSTM with Fine Tuning + Synthetic augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.185	0.175	0.195	0.118
MMI	0.081	–	0.167	0.137	0.126
OULU	0.081	0.143	–	0.137	0.126
MUG	0.269	0.208	0.167	–	0.177
AFEW	0.146	0.167	0.167	0.180	–

Table D.16. Classification accuracies of ResNet101-LSTM with Fine Tuning + Data augmentation + Synthetic augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.137	0.167	0.164	0.113
MMI	0.081	–	0.167	0.137	0.126
OULU	0.269	0.208	–	0.190	0.177
MUG	0.081	0.143	0.167	–	0.126
AFEW	0.000	0.000	0.000	0.000	–

Table D.17. Classification accuracies of C3D with Fine Tuning when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.399	0.417	0.640	0.166
MMI	0.531	–	0.192	0.473	0.173
OULU	0.906	0.167	–	0.143	0.152
MUG	0.663	0.458	0.335	–	0.178
AFEW	0.146	0.167	0.167	0.180	–

Table D.18. Classification accuracies of C3D with Fine Tuning + Data augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.310	0.460	0.502	0.158
MMI	0.460	–	0.250	0.407	0.138
OULU	0.223	0.179	–	0.187	0.137
MUG	0.269	0.208	0.167	–	0.177
AFEW	0.120	0.101	0.123	0.148	–

Table D.19. Classification accuracies of C3D with Fine Tuning + Synthetic augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.399	0.446	0.544	0.128
MMI	0.269	–	0.167	0.190	0.177
OULU	0.223	0.179	–	0.187	0.137
MUG	0.223	0.179	0.167	–	0.137
AFEW	0.146	0.167	0.167	0.180	–

Table D.20. Classification accuracies of C3D with Fine Tuning + Data augmentation + Synthetic augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.387	0.467	0.587	0.148
MMI	0.269	–	0.167	0.190	0.177
OULU	0.091	0.167	–	0.143	0.152
MUG	0.269	0.208	0.167	–	0.177
AFEW	0.146	0.167	0.167	0.180	–

Table D.21. Classification accuracies of C3D-Block-LSTM with Fine Tuning when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.440	0.410	0.616	0.161
MMI	0.563	–	0.246	0.535	0.164
OULU	0.654	0.268	–	0.566	0.166
MUG	0.634	0.482	0.404	–	0.165
AFEW	0.259	0.131	0.152	0.123	–

Table D.22. Classification accuracies of C3D-Block-LSTM with Fine Tuning + Data augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.417	0.498	0.543	0.160
MMI	0.595	–	0.271	0.481	0.148
OULU	0.725	0.327	–	0.586	0.169
MUG	0.612	0.446	0.321	–	0.140
AFEW	0.126	0.163	0.123	0.169	–

Table D.23. Classification accuracies of C3D-Block-LSTM with Fine Tuning + Synthetic augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.446	0.440	0.630	0.151
MMI	0.269	–	0.167	0.190	0.177
OULU	0.146	0.167	–	0.180	0.194
MUG	0.269	0.208	0.167	–	0.177
AFEW	0.146	0.167	0.167	0.180	–

Table D.24. Classification accuracies of C3D-Block-LSTM with Fine Tuning + Data augmentation + Synthetic augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.423	0.463	0.607	0.174
MMI	0.146	–	0.167	0.180	0.194
OULU	0.146	0.167	–	0.180	0.194
MUG	0.146	0.167	0.167	–	0.194
AFEW	0.146	0.167	0.167	0.180	–

Table D.25. Classification accuracies of I3D with Fine Tuning when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.256	0.306	0.209	0.173
MMI	0.359	–	0.269	0.327	0.163
OULU	0.505	0.345	–	0.331	0.164
MUG	0.395	0.482	0.285	–	0.165
AFEW	0.000	0.000	0.000	0.000	–

Table D.26. Classification accuracies of I3D with Fine Tuning + Data augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.280	0.294	0.272	0.173
MMI	0.366	–	0.235	0.343	0.156
OULU	0.573	0.280	–	0.371	0.140
MUG	0.330	0.440	0.277	–	0.180
AFEW	–	–	–	–	–

Table D.27. Classification accuracies of I3D with Fine Tuning + Synthetic augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.220	0.383	0.249	0.161
MMI	0.081	–	0.167	0.137	0.126
OULU	0.265	0.208	–	0.190	0.177
MUG	0.129	0.179	0.171	–	0.177
AFEW	0.000	0.000	0.000	0.000	–

Table D.28. Classification accuracies of I3D with Fine Tuning + Data augmentation + Synthetic augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.214	0.185	0.198	0.167
MMI	0.223	–	0.169	0.176	0.146
OULU	0.081	0.143	–	0.145	0.142
MUG	0.184	0.137	0.171	–	0.130
AFEW	0.000	0.000	0.000	0.000	–

Table D.29. Classification accuracies of ResNet3D-18 with Fine Tuning when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.292	0.360	0.451	0.137
MMI	0.272	–	0.210	0.336	0.181
OULU	0.443	0.387	–	0.405	0.136
MUG	0.709	0.458	0.354	–	0.199
AFEW	0.181	0.185	0.194	0.194	–

Table D.30. Classification accuracies of ResNet3D-18 with Fine Tuning + Data augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.333	0.383	0.447	0.152
MMI	0.330	–	0.281	0.366	0.154
OULU	0.631	0.268	–	0.380	0.135
MUG	0.583	0.488	0.273	–	0.124
AFEW	0.217	0.167	0.175	0.194	–

Table D.31. Classification accuracies of ResNet3D-18 with Fine Tuning + Synthetic augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.238	0.252	0.304	0.144
MMI	0.155	–	0.200	0.192	0.143
OULU	0.194	0.137	–	0.175	0.123
MUG	0.236	0.179	0.138	–	0.132
AFEW	0.168	0.155	0.165	0.168	–

Table D.32. Classification accuracies of ResNet3D-18 with Fine Tuning + Data augmentation + Synthetic augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.256	0.363	0.284	0.128
MMI	0.188	–	0.167	0.167	0.112
OULU	0.168	0.143	–	0.189	0.199
MUG	0.162	0.173	0.171	–	0.191
AFEW	0.217	0.167	0.167	0.180	–

Table D.33. Classification accuracies of ResNet3D-101 with Fine Tuning when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.393	0.460	0.519	0.157
MMI	0.417	–	0.360	0.625	0.177
OULU	0.848	0.345	–	0.424	0.131
MUG	0.586	0.506	0.275	–	0.215
AFEW	0.227	0.315	0.179	0.311	–

Table D.34. Classification accuracies of ResNet3D-101 with Fine Tuning + Data augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.411	0.433	0.535	0.180
MMI	0.476	–	0.292	0.545	0.173
OULU	0.877	0.435	–	0.527	0.156
MUG	0.696	0.548	0.358	–	0.162
AFEW	0.172	0.232	0.204	0.182	–

Table D.35. Classification accuracies of ResNet3D-101 with Fine Tuning + Synthetic augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.393	0.496	0.430	0.197
MMI	0.091	–	0.167	0.143	0.147
OULU	0.091	0.167	–	0.143	0.152
MUG	0.191	0.137	0.167	–	0.133
AFEW	0.146	0.167	0.167	0.180	–

Table D.36. Classification accuracies of ResNet3D-101 with Fine Tuning + Data augmentation + Synthetic augmentation when trained in one database and evaluated in others.

Train	CK+	MMI	OULU	MUG	AFEW
CK+	–	0.310	0.550	0.512	0.186
MMI	0.223	–	0.167	0.187	0.138
OULU	0.091	0.167	–	0.143	0.152
MUG	0.220	0.167	0.160	–	0.155
AFEW	0.146	0.167	0.167	0.180	–

Table E.3. p -values of the hypothesis testing for (a) $H_0 : M_1 = M_2$, $H_a : M_1 \neq M_2$ and (b) $H_0 : M_1 \leq M_2$, $H_a : M_1 > M_2$ on CK+ for single-vs-merged-database accuracies, where M_1 are the methods on the left and M_2 are the methods on top. The shaded cells denote the rejection of the null hypothesis at 5% significance level.

(a) $H_0 : M_1 = M_2$, $H_a : M_1 \neq M_2$

	VGGrLSTM										InceptionV3-LSTM										ResNet18-LSTM										ResNet101-LSTM										CID										C3D-Block-LSTM										iD										ResNet101-18										ResNet101-101									
	RI	RI+D	RI+SD	FT	FT+DA	FT+SD	RI	RI+D	RI+SD	FT	FT+DA	FT+SD	RI	RI+D	RI+SD	FT	FT+DA	FT+SD	RI	RI+D	RI+SD	FT	FT+DA	FT+SD	RI	RI+D	RI+SD	FT	FT+DA	FT+SD	RI	RI+D	RI+SD	FT	FT+DA	FT+SD	RI	RI+D	RI+SD	FT	FT+DA	FT+SD	RI	RI+D	RI+SD	FT	FT+DA	FT+SD	RI	RI+D	RI+SD	FT	FT+DA	FT+SD	RI	RI+D	RI+SD	FT	FT+DA	FT+SD																														
VGGrLSTM	RI	0.04	0.04	0.04	0.04	0.04	RI+D	0.04	0.04	0.04	0.04	RI+SD	0.04	0.04	0.04	0.04	FT	0.04	0.04	0.04	0.04	0.04	0.04	FT+DA	0.04	0.04	0.04	0.04	0.04	FT+SD	0.04	0.04	0.04	0.04	0.04	0.04	RI	0.04	0.04	0.04	0.04	0.04	0.04	RI+D	0.04	0.04	0.04	0.04	0.04	RI+SD	0.04	0.04	0.04	0.04	0.04	FT	0.04	0.04	0.04	0.04	0.04	0.04	FT+DA	0.04	0.04	0.04	0.04	0.04	FT+SD	0.04	0.04	0.04	0.04	0.04	0.04															

(b) $H_0 : M_1 \leq M_2$, $H_a : M_1 > M_2$

	VGGrLSTM										InceptionV3-LSTM										ResNet18-LSTM										ResNet101-LSTM										CID										C3D-Block-LSTM										iD										ResNet101-18										ResNet101-101									
	RI	RI+D	RI+SD	FT	FT+DA	FT+SD	RI	RI+D	RI+SD	FT	FT+DA	FT+SD	RI	RI+D	RI+SD	FT	FT+DA	FT+SD	RI	RI+D	RI+SD	FT	FT+DA	FT+SD	RI	RI+D	RI+SD	FT	FT+DA	FT+SD	RI	RI+D	RI+SD	FT	FT+DA	FT+SD	RI	RI+D	RI+SD	FT	FT+DA	FT+SD	RI	RI+D	RI+SD	FT	FT+DA	FT+SD	RI	RI+D	RI+SD	FT	FT+DA	FT+SD	RI	RI+D	RI+SD	FT	FT+DA	FT+SD																														
VGGrLSTM	RI	0.02	0.02	0.02	0.02	0.02	RI+D	0.02	0.02	0.02	0.02	RI+SD	0.02	0.02	0.02	0.02	FT	0.02	0.02	0.02	0.02	0.02	0.02	FT+DA	0.02	0.02	0.02	0.02	0.02	FT+SD	0.02	0.02	0.02	0.02	0.02	0.02	RI	0.02	0.02	0.02	0.02	0.02	0.02	RI+D	0.02	0.02	0.02	0.02	0.02	RI+SD	0.02	0.02	0.02	0.02	0.02	FT	0.02	0.02	0.02	0.02	0.02	0.02	FT+DA	0.02	0.02	0.02	0.02	0.02	FT+SD	0.02	0.02	0.02	0.02	0.02	0.02															

Table E.13. p -values of the hypothesis testing for (a) $H_0 : M_1 = M_2, H_a : M_1 \neq M_2$ and (b) $H_0 : M_1 \leq M_2, H_a : M_1 > M_2$ for single-models-unified-database accuracies, where M_1 are the methods on the left and M_2 are the methods on top. The shaded cells denote the rejection of the null hypothesis at 5% significance level.

(a) $H_0 : M_1 = M_2, H_a : M_1 \neq M_2$

		CK+			MMI						OULU-Casia						MUG									
		RI	RI+DA	RI+SD	FT	FT+DA	FT+SD	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD	
CK+	RI	-	0.15	0.10	0.00	0.00	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
	RI+DA	0.15	-	0.39	0.00	0.00	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
	RI+SD	0.10	0.39	-	0.00	0.00	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	FT	0.00	0.00	0.00	-	0.13	0.07	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	FT+DA	0.00	0.00	0.00	0.13	-	0.76	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	FT+SD	0.00	0.00	0.00	0.07	0.76	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
MMI	RI	-	-	-	-	-	-	0.07	0.91	0.00	0.00	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-	
	RI+DA	-	-	-	-	-	0.07	-	0.16	0.00	0.00	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-	
	RI+SD	-	-	-	-	-	0.91	0.16	-	0.00	0.00	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-	
	FT	-	-	-	-	-	0.00	0.00	0.00	-	0.73	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-	
	FT+DA	-	-	-	-	-	0.00	0.00	0.00	0.73	-	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-	
	FT+SD	-	-	-	-	-	0.00	0.00	0.00	0.00	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
OULU-Casia	RI	-	-	-	-	-	-	-	-	-	-	-	-	0.14	0.60	0.00	0.00	0.00	-	-	-	-	-	-	-	
	RI+DA	-	-	-	-	-	-	-	-	-	-	-	0.14	-	0.62	0.00	0.00	0.00	-	-	-	-	-	-	-	
	RI+SD	-	-	-	-	-	-	-	-	-	-	-	0.60	0.62	-	0.00	0.00	0.00	-	-	-	-	-	-	-	
	FT	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	-	0.79	0.99	-	-	-	-	-	-	
	FT+DA	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.79	-	0.16	-	-	-	-	-	-	
	FT+SD	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.99	0.16	-	-	-	-	-	-	-	
MUG	RI	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.93	0.55	0.00	0.00	0.00	
	RI+DA	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.93	-	0.36	0.00	0.00	
	RI+SD	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.55	0.36	-	0.00	0.00	
	FT	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	-	0.00	
	FT+DA	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.00	-	
	FT+SD	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.01	0.82	

(b) $H_0 : M_1 \leq M_2, H_a : M_1 > M_2$

		CK+			MMI						OULU-Casia						MUG								
		RI	RI+DA	RI+SD	FT	FT+DA	FT+SD	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD
CK+	RI	-	0.08	0.05	1.00	1.00	1.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	RI+DA	0.92	-	0.20	1.00	1.00	1.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	RI+SD	0.95	0.80	-	1.00	1.00	1.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	FT	0.00	0.00	0.00	-	0.93	0.96	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	FT+DA	0.00	0.00	0.00	0.07	-	0.62	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	FT+SD	0.00	0.00	0.00	0.04	0.38	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
MMI	RI	-	-	-	-	-	-	0.96	0.54	1.00	1.00	1.00	-	-	-	-	-	-	-	-	-	-	-	-	-
	RI+DA	-	-	-	-	-	0.04	-	0.08	1.00	1.00	1.00	-	-	-	-	-	-	-	-	-	-	-	-	-
	RI+SD	-	-	-	-	-	0.46	0.92	-	1.00	1.00	1.00	-	-	-	-	-	-	-	-	-	-	-	-	-
	FT	-	-	-	-	-	0.00	0.00	0.00	-	0.64	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-
	FT+DA	-	-	-	-	-	0.00	0.00	0.00	0.36	-	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-
	FT+SD	-	-	-	-	-	0.00	0.00	0.00	1.00	1.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-
OULU-Casia	RI	-	-	-	-	-	-	-	-	-	-	-	-	0.07	0.30	1.00	1.00	1.00	-	-	-	-	-	-	-
	RI+DA	-	-	-	-	-	-	-	-	-	-	-	0.93	-	0.69	1.00	1.00	1.00	-	-	-	-	-	-	-
	RI+SD	-	-	-	-	-	-	-	-	-	-	-	0.70	0.31	-	1.00	1.00	1.00	-	-	-	-	-	-	-
	FT	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	-	0.40	0.51	-	-	-	-	-	-
	FT+DA	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.60	-	0.92	-	-	-	-	-	-
	FT+SD	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.49	0.08	-	-	-	-	-	-	-
MUG	RI	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.46	0.72	1.00	1.00	1.00
	RI+DA	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.54	-	0.82	1.00	1.00
	RI+SD	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.28	0.18	-	1.00	1.00
	FT	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	-	0.00
	FT+DA	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	1.00	-
	FT+SD	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.99	0.59

Table E.14. p -values of the hypothesis testing for (a) $H_0 : M_1 = M_2, H_a : M_1 \neq M_2$ and (b) $H_0 : M_1 \leq M_2, H_a : M_1 > M_2$ for merged-models-unified-database accuracies, where M_1 are the methods on the left and M_2 are the methods on top. The shaded cells denote the rejection of the null hypothesis at 5% significance level.

(a) $H_0 : M_1 = M_2, H_a : M_1 \neq M_2$

		CK+						MMI						OULU-Casia						MUG					
		RI	RI+DA	RI+SD	FT	FT+DA	FT+SD	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD
CK+	RI	-	0.03	0.10	0.00	0.00	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	RI+DA	0.03	-	0.10	0.00	0.00	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	RI+SD	0.10	0.10	-	0.00	0.00	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	FT	0.00	0.00	0.00	-	0.02	0.67	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	FT+DA	0.00	0.00	0.00	0.02	-	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	FT+SD	0.00	0.00	0.00	0.67	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
MMI	RI	-	-	-	-	-	-	0.00	0.55	0.00	0.35	0.00	-	-	-	-	-	-	-	-	-	-	-	-	
	RI+DA	-	-	-	-	-	-	0.00	0.00	0.00	0.17	0.00	-	-	-	-	-	-	-	-	-	-	-	-	
	RI+SD	-	-	-	-	-	-	0.55	0.00	-	0.00	0.23	0.00	-	-	-	-	-	-	-	-	-	-	-	-
	FT	-	-	-	-	-	-	0.00	0.00	0.00	-	0.00	0.68	-	-	-	-	-	-	-	-	-	-	-	-
	FT+DA	-	-	-	-	-	-	0.35	0.17	0.23	0.00	-	0.00	-	-	-	-	-	-	-	-	-	-	-	-
	FT+SD	-	-	-	-	-	-	0.00	0.00	0.00	0.68	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-
OULU-Casia	RI	-	-	-	-	-	-	-	-	-	-	-	0.00	0.17	0.00	0.02	0.00	-	-	-	-	-	-		
	RI+DA	-	-	-	-	-	-	-	-	-	-	-	0.00	-	0.00	0.00	0.01	0.00	-	-	-	-	-	-	
	RI+SD	-	-	-	-	-	-	-	-	-	-	-	0.17	0.00	-	0.00	0.07	0.00	-	-	-	-	-	-	
	FT	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	-	0.00	0.14	-	-	-	-	-	-	
	FT+DA	-	-	-	-	-	-	-	-	-	-	-	0.02	0.01	0.07	0.00	-	0.00	-	-	-	-	-	-	
	FT+SD	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.14	0.00	-	-	-	-	-	-	-	
MUG	RI	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.43	0.00	0.00	0.00	0.00		
	RI+DA	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	-	0.00	0.00	0.02	0.00		
	RI+SD	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.43	0.00	-	0.00	0.00	0.00		
	FT	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	-	0.00	0.71		
	FT+DA	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.02	0.00	0.00	-	0.00		
	FT+SD	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.71	0.00	-		

(b) $H_0 : M_1 \leq M_2, H_a : M_1 > M_2$

		CK+						MMI						OULU-Casia						MUG					
		RI	RI+DA	RI+SD	FT	FT+DA	FT+SD	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD
CK+	RI	-	0.01	0.05	1.00	1.00	1.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	RI+DA	0.99	-	0.95	1.00	1.00	1.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	RI+SD	0.95	0.05	-	1.00	1.00	1.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	FT	0.00	0.00	0.00	-	0.01	0.67	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	FT+DA	0.00	0.00	0.00	0.99	-	1.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
	FT+SD	0.00	0.00	0.00	0.33	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
MMI	RI	-	-	-	-	-	-	0.00	0.73	1.00	0.18	1.00	-	-	-	-	-	-	-	-	-	-	-	-	
	RI+DA	-	-	-	-	-	-	1.00	-	1.00	1.00	0.92	1.00	-	-	-	-	-	-	-	-	-	-	-	-
	RI+SD	-	-	-	-	-	-	0.27	0.00	-	1.00	0.11	1.00	-	-	-	-	-	-	-	-	-	-	-	-
	FT	-	-	-	-	-	-	0.00	0.00	0.00	-	0.00	0.66	-	-	-	-	-	-	-	-	-	-	-	-
	FT+DA	-	-	-	-	-	-	0.82	0.08	0.89	1.00	-	1.00	-	-	-	-	-	-	-	-	-	-	-	-
	FT+SD	-	-	-	-	-	-	0.00	0.00	0.00	0.34	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-
OULU-Casia	RI	-	-	-	-	-	-	-	-	-	-	-	0.00	0.09	1.00	0.01	1.00	-	-	-	-	-	-		
	RI+DA	-	-	-	-	-	-	-	-	-	-	-	1.00	-	1.00	1.00	1.00	1.00	-	-	-	-	-	-	
	RI+SD	-	-	-	-	-	-	-	-	-	-	-	0.91	0.00	-	1.00	0.04	1.00	-	-	-	-	-	-	
	FT	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	-	0.00	0.93	-	-	-	-	-	-	
	FT+DA	-	-	-	-	-	-	-	-	-	-	-	0.99	0.00	0.96	1.00	-	1.00	-	-	-	-	-	-	
	FT+SD	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.07	0.00	-	-	-	-	-	-	-	
MUG	RI	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.22	1.00	0.00	1.00	1.00		
	RI+DA	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1.00	-	1.00	1.00	0.99	1.00		
	RI+SD	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.78	0.00	-	1.00	0.00	1.00		
	FT	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	-	0.00	0.64		
	FT+DA	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1.00	0.01	1.00	1.00	-	1.00		
	FT+SD	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.36	0.00	-		

Table E.15. p -values of the hypothesis testing for (a) $H_0 : M_1 = M_2$, $H_a : M_1 \neq M_2$ and (b) $H_0 : M_1 \leq M_2$, $H_a : M_1 > M_2$ for single-vs-merged-models-unified-database accuracies, where M_1 are the methods on the left and M_2 are the methods on top. The shaded cells denote the rejection of the null hypothesis at 5% significance level.

(a) $H_0 : M_1 = M_2$, $H_a : M_1 \neq M_2$

		CK+						MMI						OULU-Casia						MUG					
		RI	RI+DA	RI+SD	FT	FT+DA	FT+SD	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD
CK+	RI	0.00	0.00	0.00	0.00	0.01	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
	RI+DA	0.01	0.00	0.00	0.00	0.00	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
	RI+SD	0.16	0.03	0.04	0.00	0.00	0.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
	FT	0.00	0.00	0.00	0.10	0.45	0.03	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
	FT+DA	0.00	0.00	0.00	0.75	0.08	0.65	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
	FT+SD	0.00	0.00	0.00	0.75	0.13	0.41	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
MMI	RI	-	-	-	-	-	-	0.02	0.18	0.01	0.00	0.83	0.00	-	-	-	-	-	-	-	-	-	-	-	
	RI+DA	-	-	-	-	-	-	0.22	0.01	0.17	0.00	0.49	0.00	-	-	-	-	-	-	-	-	-	-	-	
	RI+SD	-	-	-	-	-	-	0.00	0.14	0.01	0.00	0.74	0.00	-	-	-	-	-	-	-	-	-	-	-	
	FT	-	-	-	-	-	-	0.00	0.00	0.00	0.00	0.00	0.00	-	-	-	-	-	-	-	-	-	-	-	
	FT+DA	-	-	-	-	-	-	0.00	0.00	0.00	0.00	0.00	0.00	-	-	-	-	-	-	-	-	-	-	-	
	FT+SD	-	-	-	-	-	-	0.00	0.00	0.00	0.00	0.00	0.00	-	-	-	-	-	-	-	-	-	-	-	
OULU-Casia	RI	-	-	-	-	-	-	-	-	-	-	-	-	0.20	0.00	0.04	0.00	0.00	0.00	-	-	-	-	-	
	RI+DA	-	-	-	-	-	-	-	-	-	-	-	-	0.34	0.00	0.01	0.00	0.01	0.00	-	-	-	-	-	
	RI+SD	-	-	-	-	-	-	-	-	-	-	-	-	0.28	0.00	0.02	0.00	0.01	0.00	-	-	-	-	-	
	FT	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.24	0.00	0.11	-	-	-	-	-	
	FT+DA	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.42	0.00	0.07	-	-	-	-	-	
	FT+SD	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.27	0.00	0.12	-	-	-	-	-	
MUG	RI	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.01	0.00	0.01	0.00	0.00	
	RI+DA	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.03	0.00	0.02	0.00	0.00	
	RI+SD	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.00	0.00	
	FT	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.91	0.00	
	FT+DA	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.05	0.00	
	FT+SD	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.02	0.00	

(b) $H_0 : M_1 \leq M_2$, $H_a : M_1 > M_2$

		CK+						MMI						OULU-Casia						MUG					
		RI	RI+DA	RI+SD	FT	FT+DA	FT+SD	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD
CK+	RI	0.00	0.00	0.00	1.00	0.99	1.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
	RI+DA	0.00	0.00	0.00	1.00	1.00	1.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
	RI+SD	0.08	0.01	0.02	1.00	1.00	1.00	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
	FT	0.00	0.00	0.00	0.95	0.23	0.99	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
	FT+DA	0.00	0.00	0.00	0.63	0.04	0.68	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
	FT+SD	0.00	0.00	0.00	0.63	0.07	0.79	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
MMI	RI	-	-	-	-	-	-	0.99	0.09	0.99	1.00	0.58	1.00	-	-	-	-	-	-	-	-	-	-	-	
	RI+DA	-	-	-	-	-	-	0.89	0.01	0.91	1.00	0.24	1.00	-	-	-	-	-	-	-	-	-	-	-	
	RI+SD	-	-	-	-	-	-	1.00	0.07	1.00	1.00	0.63	1.00	-	-	-	-	-	-	-	-	-	-	-	
	FT	-	-	-	-	-	-	0.00	0.00	0.00	1.00	0.00	1.00	-	-	-	-	-	-	-	-	-	-	-	
	FT+DA	-	-	-	-	-	-	0.00	0.00	0.00	1.00	0.00	1.00	-	-	-	-	-	-	-	-	-	-	-	
	FT+SD	-	-	-	-	-	-	0.00	0.00	0.00	1.00	0.00	1.00	-	-	-	-	-	-	-	-	-	-	-	
OULU-Casia	RI	-	-	-	-	-	-	-	-	-	-	-	-	0.10	0.00	0.02	1.00	0.00	1.00	-	-	-	-	-	
	RI+DA	-	-	-	-	-	-	-	-	-	-	-	-	0.17	0.00	0.01	1.00	0.01	1.00	-	-	-	-	-	
	RI+SD	-	-	-	-	-	-	-	-	-	-	-	-	0.14	0.00	0.01	1.00	0.00	1.00	-	-	-	-	-	
	FT	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.88	0.00	0.95	-	-	-	-	-	
	FT+DA	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.79	0.00	0.97	-	-	-	-	-	
	FT+SD	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.87	0.00	0.94	-	-	-	-	-	
MUG	RI	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.01	0.00	0.00	1.00	0.00	
	RI+DA	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.02	0.00	0.01	1.00	0.00	
	RI+SD	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	1.00	0.00	
	FT	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.45	0.00	
	FT+DA	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.97	0.00	
	FT+SD	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0.00	0.00	0.00	0.99	0.00	

Table E.16. p -values of the hypothesis testing for (a) $H_0 : M_1 = M_2, H_a : M_1 \neq M_2$ and (b) $H_0 : M_1 \leq M_2, H_a : M_1 > M_2$ on single-models unified-dataset accuracies, where M_1 are the methods on the left and M_2 are the methods on top. The shaded cells denote the rejection of the null hypothesis at 5% significance level.

(a) $H_0 : M_1 = M_2, H_a : M_1 \neq M_2$						
	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD
RI	–	0.37	0.44	0.00	0.00	0.00
RI+DA	0.37	–	0.90	0.00	0.00	0.00
RI+SD	0.44	0.90	–	0.00	0.00	0.00
FT	0.00	0.00	0.00	–	0.67	0.05
FT+DA	0.00	0.00	0.00	0.67	–	0.27
FT+SD	0.00	0.00	0.00	0.05	0.27	–
(b) $H_0 : M_1 \leq M_2, H_a : M_1 > M_2$						
	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD
RI	–	0.18	0.22	1.00	1.00	1.00
RI+DA	0.82	–	0.45	1.00	1.00	1.00
RI+SD	0.78	0.55	–	1.00	1.00	1.00
FT	0.00	0.00	0.00	–	0.34	0.02
FT+DA	0.00	0.00	0.00	0.66	–	0.13
FT+SD	0.00	0.00	0.00	0.98	0.87	–

Table E.17. p -values of the hypothesis testing for (a) $H_0 : M_1 = M_2, H_a : M_1 \neq M_2$ and (b) $H_0 : M_1 \leq M_2, H_a : M_1 > M_2$ on merged-models unified-dataset accuracies, where M_1 are the methods on the left and M_2 are the methods on top. The shaded cells denote the rejection of the null hypothesis at 5% significance level.

(a) $H_0 : M_1 = M_2, H_a : M_1 \neq M_2$						
	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD
RI	–	0.00	0.14	0.00	0.81	0.00
RI+DA	0.00	–	0.00	0.00	0.00	0.00
RI+SD	0.14	0.00	–	0.00	0.91	0.00
FT	0.00	0.00	0.00	–	0.00	0.24
FT+DA	0.81	0.00	0.91	0.00	–	0.00
FT+SD	0.00	0.00	0.00	0.24	0.00	–
(b) $H_0 : M_1 \leq M_2, H_a : M_1 > M_2$						
	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD
RI	–	0.00	0.07	1.00	0.40	1.00
RI+DA	1.00	–	1.00	1.00	1.00	1.00
RI+SD	0.93	0.00	–	1.00	0.55	1.00
FT	0.00	0.00	0.00	–	0.00	0.88
FT+DA	0.60	0.00	0.45	1.00	–	1.00
FT+SD	0.00	0.00	0.00	0.12	0.00	–

Table E.18. p -values of the hypothesis testing for (a) $H_0 : M_1 = M_2$, $H_a : M_1 \neq M_2$ and (b) $H_0 : M_1 \leq M_2$, $H_a : M_1 > M_2$ on single-vs-merged unified-models and -dataset accuracies, where M_1 are the methods on the left and M_2 are the methods on top. The shaded cells denote the rejection of the null hypothesis at 5% significance level.

(a) $H_0 : M_1 = M_2$, $H_a : M_1 \neq M_2$						
	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD
RI	0.00	0.00	0.00	0.00	0.03	0.00
RI+DA	0.01	0.00	0.00	0.00	0.11	0.00
RI+SD	0.07	0.00	0.00	0.00	0.12	0.00
FT	0.00	0.00	0.00	0.00	0.00	0.00
FT+DA	0.00	0.00	0.00	0.00	0.00	0.00
FT+SD	0.00	0.00	0.00	0.00	0.00	0.00
(b) $H_0 : M_1 \leq M_2$, $H_a : M_1 > M_2$						
	RI	RI+DA	RI+SD	FT	FT+DA	FT+SD
RI	0.00	0.00	0.00	1.00	0.01	1.00
RI+DA	0.01	0.00	0.00	1.00	0.05	1.00
RI+SD	0.03	0.00	0.00	1.00	0.06	1.00
FT	0.00	0.00	0.00	1.00	0.00	1.00
FT+DA	0.00	0.00	0.00	1.00	0.00	1.00
FT+SD	0.00	0.00	0.00	1.00	0.00	1.00