

The Effect of Model Compression on Fairness in Facial Expression Recognition

Samuil Stoychev
University of Cambridge
Cambridge, United Kingdom
ss2719@cantab.ac.uk

Hatice Gunes
University of Cambridge
Cambridge, United Kingdom
hatice.gunes@cl.cam.ac.uk

ABSTRACT

Deep neural networks have proved hugely successful, achieving human-like performance on a variety of tasks. However, they are also computationally expensive, which has motivated the development of model compression techniques which reduce the resource consumption associated with deep learning models. Nevertheless, recent studies have suggested that model compression can have an adverse effect on algorithmic fairness, amplifying existing biases in machine learning models. With this project we aim to extend those studies to the context of facial expression recognition. To do that, we set up a neural network classifier to perform facial expression recognition and implement several model compression techniques on top of it. We then run experiments on two facial expression datasets, namely the Extended Cohn-Kanade Dataset (CK+DB) and the Real-World Affective Faces Database (RAF-DB), to examine the individual and combined effect that compression techniques have on the model size, accuracy and fairness. Our experimental results show that: (i) Compression and quantisation achieve significant reduction in model size with minimal impact on overall accuracy for both CK+DB and RAF-DB; (ii) in terms of model accuracy, the classifier trained and tested on RAF-DB seems more robust to compression compared to the CK+ DB; (iii) for RAF-DB, the different compression strategies do not seem to increase the gap in predictive performance across the sensitive attributes of gender, race and age which is in contrast with the results on the CK+DB, where compression seems to amplify existing biases for gender. We analyse the results and discuss the potential reasons for our findings.

CCS CONCEPTS

• **Computing methodologies** → **Neural networks**; • **Human-centered computing**;

KEYWORDS

facial expression recognition, algorithmic fairness, neural networks, model compression

1 INTRODUCTION

Recent years have seen *deep neural networks* (DNNs) achieve state-of-the-art performance on a variety of problems including face recognition [38], cancer detection [27], natural language processing [27], etc. Deep learning has proved particularly effective at extracting meaningful representations from raw data [42].

However, as the predictive performance of deep neural networks has increased, so has the size of deep learning architectures: Modern DNNs can consist of hundreds of millions of parameters [11],

making them slow to train and hard to store. Deep learning’s growing computational cost has made it hard to deploy deep learning models on resource-constrained devices (e.g. mobile phones, robots, microcontrollers) which often lack the storage, memory or processing power to support large DNNs [23, 28, 45]. The high resource consumption associated with deep learning models has also been problematic in the light of initiatives such as the “Green-AI” [41, 46] movement advocating for a reduction in the carbon emissions and the environmental impact associated with artificial intelligence.

This has given rise to the development of *model compression* strategies, which aim to reduce the size of deep learning models. Examples of model compression techniques include *pruning* [57], *quantisation* [23], *weight clustering* [20], etc. We provide a more detailed overview of the compression strategies considered in this project in Section 2.2.

However, a couple of recent studies have suggested that, by reducing the network capacity of the DNNs, model compression can amplify existing biases: Hooker et al. [22] demonstrate that pruning and post-training quantisation can amplify biases when classifying hair colour on CelebA. This issue is also raised in a study by Paganini [37] who discusses the effect of pruning on algorithmic fairness and proposes a framework for fair model pruning.

With this work, we aim to extend the aforementioned studies by Paganini and Hooker et al. to the context of affective computing. In particular, we consider the task of *facial expression recognition* (FER) where the model has to classify expressions based on images of human faces. To this end, we train FER models using the CK+DB and the RAF-DB, and implement three compression strategies (pruning, weight clustering and post-training quantisation) on top of them. We then evaluate and compare the performance of the baseline models against the performance of the compressed models, and analyse the results to address three research questions:

- **RQ1: “How effective is model compression in the context of FER?”** That is, can compression techniques achieve a considerable reduction in the model size, while preserving a high level of predictive accuracy?
- **RQ2: “Do model compression techniques amplify biases?”** Here we seek to verify the claims by Paganini and Hooker et al. across a wider variety of compression techniques and in the context of FER.
- **RQ3: “Is the impact on fairness identical across different compression techniques?”** We are interested to know whether all compression strategies amplify biases to the same extent.

This study extends the previous works by Paganini and Hooker et al. in three directions:

- **Extending the problem to affective computing:** We consider the problem of compression’s effect on fairness in the context of affective computing and, in particular, facial expression recognition. By comparison, the study by Hooker et al. is based on classifying hair colour on CelebA [33], and the study by Paganini considers object recognition and digit classification tasks.

The topic of model compression is particularly relevant to the field of affective computing as affective computation is increasingly performed on constrained devices such as robots [43] and mobile phones [19].

- **Considering more model compression techniques:** Our study involves three compression strategies (pruning, weight clustering and post-training quantisation) – one more by Hooker et al. (who consider pruning and post-training quantisation), and two more by the study by Paganini, which focuses solely on pruning.
- **Considering the combined effect of compression techniques:** In practice, compression techniques are often combined together to form so called “*compression pipelines*” [20]. That is why, we also consider two combinations of compression strategies (pruning with quantisation and weight clustering with quantisation). The previous studies have only examined the individual behaviour of compression techniques.

2 BACKGROUND

2.1 Algorithmic Fairness

Nowadays, machine learning algorithms are used to inform or automate decision-making across various fields of high social importance. Machine learning approaches have been used for automating recruitment in large companies [53], assigning credit scores [30, 54] and anticipating criminal activity [21] to name a few.

The increasing impact of machine learning on our society has highlighted the importance of *algorithmic fairness*. An algorithm is considered to be fair if its behaviour is not improperly influenced by sensitive attributes (e.g., a person’s gender or race) [36].

Nevertheless, recent studies have exposed the propensity of algorithms to be unfair and exhibit dangerous *biases*, potentially “*reinforcing the discriminatory practices in society*” [6]. For example, Amazon’s AI recruitment tool has been reported to favour male applicants over female applicants [15]. Apple’s credit score has also been shown to systematically disadvantage women [1]. A study by Joy Buolamwini has demonstrated that popular facial analysis services perform disproportionately poorly on dark-skinned females [10].

The increasing awareness of algorithmic biases has given rise to multiple fairness initiatives such the Algorithmic Justice League [24], IBM’s AI Fairness 360 [2] and Google’s ML Fairness [44]. Research into fairness in facial expression recognition has also started gaining momentum. Xu et al. [50] compared three different bias-mitigation approaches, namely, a baseline, an attribute-aware, and a disentangled approach on two well-known data sets, Real-World Affected Faces-Database and CelebFaces Attributes. Cheong et al. in [13] provided an overview and techniques that can be used for achieving fairness in facial affect recognition.

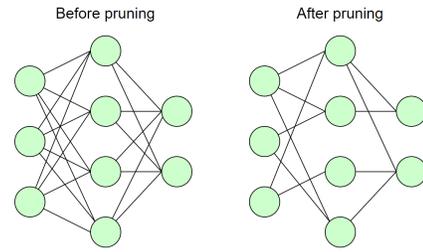


Figure 1: A DNN before and after pruning has been applied. The network has been stripped down to a subset of its original weights and the overall density has decreased. Illustration adapted from [47].

Despite the large body of research which has studied the problem, though, there is still no consensus in the scientific community on what the precise definition of fairness should be. Multiple definitions of fairness have been proposed but none of them is a “silver bullet” that fits all use cases. Instead, the “right” choice of a fairness metric often depends on the specific context in which the algorithm is used [7].

For this work, and in the context of facial expression recognition, we adopt the fairness definition of *overall accuracy equality* [48]. The definition is akin to the concepts of predictive parity [39] and disparate mistreatment [55], and states that *a fair algorithm should have the same predictive accuracy regardless of any underlying sensitive attributes*.

To express this formally, assume we have a facial expression recognition model that aims to predict a subject’s true expression Y by producing a prediction \hat{Y} . Let the subject’s gender be denoted by G and be equal to either m (male) or f (female)¹. In that case, we expect that a fair model would have the following property:

$$P(Y = \hat{Y} | G = m) = P(Y = \hat{Y} | G = f) \quad (1)$$

That is, we expect that a fair FER model would classify the expressions of male and female subjects with the same accuracy.

2.2 Model Compression Techniques

2.2.1 Quantisation. Quantisation is a popular compression method which can significantly reduce the size of the model, leading to savings in storage and memory [12]. The key idea behind quantisation is sacrificing precision for efficiency – while most standard DNN implementations represent weights and activations using the `float32` datatype, quantisation allows to represent those values using a smaller data type – normally `float16` or `int8` [32]. Quantisation can either be introduced *during* training (also known as *quantisation-aware training*), or it can be applied to a pre-trained model, which is known as *post-training quantisation* [35]. In our experiments, we consider post-training quantisation to the `int8` type.

¹Scheuerman et al. [40] have noted that in the context of facial analysis, it is useful to differentiate *gender appearance* (whether a subject appears feminine or masculine) from *gender self-identification* (whether a subject identifies as male or female). Throughout this report, we use the terms “male” and “female” to refer only to gender appearance and make no assumptions about the gender self-identification of the subjects.

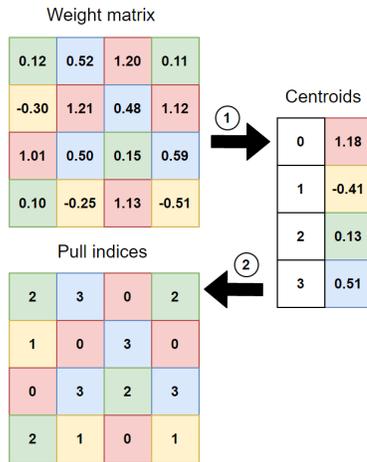


Figure 2: Diagram depicting the process of weight clustering. Similar weights get grouped together into the same cluster (step ①), after which weights are replaced with pull indices (step ②). Figure adapted from [3]

2.2.2 Pruning. Weight pruning is another compression strategy which can greatly reduce the size of a DNN [32]. It does so by eliminating redundant weights which contribute little to the behaviour of the model. As illustrated in Figure 1, pruning reduces the density of the neural network, making it more lightweight and easier to compress by traditional compression tools such as zip². Which weights get pruned is dictated by the pruning strategy. The most popular approach, which we focus on in this project, is called *magnitude-based pruning* [9] and eliminates the weights with the lowest absolute value – i.e., the ones whose values are the closest to zero. The proportion of weights that need to be pruned is called the pruning *sparsity* – for example, pruning at 90% sparsity would remove 90% of the connections in a given network.

2.2.3 Weight Clustering. Weight clustering (also known as *weight sharing* [20]) reduces the size of the model by grouping together weight of similar values. This process is illustrated in Figure 2: During step ①, weight matrices are processed by a clustering algorithm, which maps each weight to one of n clusters (where n is the number of clusters specified by the user). Each cluster consists of an *index* (one of 0, 1, ..., $n - 1$) and a *centroid* value which is representative of the values of the weights mapped to that cluster. During step ②, the weight matrices of the DNN are replaced with *pull indices* – instead of containing the values of the corresponding weights, each pull index contains the index of the cluster containing the respective weight. During inference, the DNN model can use the pull indices to obtain the centroid values corresponding to each weight.

Clustering reduces the model size for two reasons: First, float values only need to be stored to represent the n centroid values. Meanwhile, the entire weight matrix is replaced with pull indices, each index represented by the smaller integer type. And second,

²https://www.tensorflow.org/model_optimization/guide/pruning

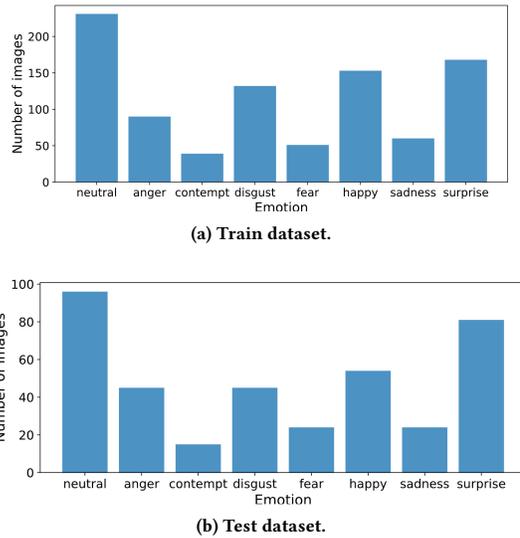


Figure 3: Distribution of emotions across the train and test split of the CK+ dataset.



Figure 4: Images from the CK+ [34] dataset after cropping.

the resulting pull indices are more likely to contain repeating values, making standard compression tools (e.g., zip) more effective, similar to pruning.

3 IMPLEMENTATION

In this chapter, we summarise the implementation steps which the project has involved. A more detailed description of the technical implementation is available in the README.md file of the code repository provided with the submission.

3.1 Data

To perform facial expression recognition, we make use of two popular datasets of human faces - CK+ [34] and RAF-DB [31], which we describe below.

3.1.1 Extended Cohn-Kanade Dataset. The Extended Cohn-Kanade Dataset (CK+) has been widely used in the context of facial expression recognition [8, 52]. It contains 327 labelled image sequences across 123 unique subjects, expressing one of 8 basic emotions - “neutral”, “anger”, “contempt”, “disgust”, “fear”, “happy”, “sadness”



Figure 5: Images from the RAF-DB [31] dataset.

and “surprise”. Images have been obtained in a controlled lab environment with subjects facing the camera as illustrated in Figure 4.

In order to examine bias, we need annotations of demographic attributes. CK+ does not provide any annotations in that respect, so we manually annotate all 123 subjects based on their gender appearance. We assign a value “male” if the subject looks masculine, and a value “female” if they look feminine. The annotations are provided in the project repository under `ckplus_labels.csv`. According to our annotation, the dataset consists of 84 female subjects and 39 male subjects.

We then apply several pre-processing steps to the CK+ dataset: For each sequence in the dataset, we take the first frame to represent a neutral emotion, and the last 3 frames to represent the emotion which the sequence was annotated with (e.g. “happy”, “sad”, etc.). That is a common pre-processing step since CK+ sequences “are from the neutral face to the peak expression” according to the dataset’s documentation. We then use the `dlib`³ library to detect the faces of the subjects and crop the images around them (allowing an extra 10% on each side to avoid cropping out parts of the chin, forehead or ears).

For validation purposes, we split the original CK+ dataset into a *train* and *test* dataset. We use cross-subject validation, allocating 86 subjects to the train dataset and the other 37 to the test dataset. That gives us 924 images in total in the train dataset and 384 images in the test dataset. The train and test dataset follow a similar distribution with respect to the emotion labels as shown in Figure 3.

Finally, several data transformations are applied to the images using TensorFlow’s `ImageDataGenerator`⁴: Images are scaled down to 48×48 pixels and converted to grayscale (since CK+ contains some RGB images). Finally, to compensate for the relatively small size of the dataset, we augment the data by applying random horizontal flipping and random rotation in the range $[-10^\circ, 10^\circ]$.

3.1.2 Real-World Affective Faces Database. The Real-World Affective Faces Database (RAF-DB) [31] is another dataset of human faces which has been commonly used in the field of FER [49, 51]. It provides 15,339 RGB images of human faces, aligned into squares of 100×100 pixels. Each image depicts one of seven basic emotions: “surprise”, “fear”, “disgust”, “happiness”, “sadness”, “anger” and “neutral”.

Unlike CK+ images, RAF-DB images are “in-the-wild” - they have not been recorded in a controlled environment, so emotions are

³<http://dlib.net/>

⁴https://www.tensorflow.org/api_docs/python/tf/keras/preprocessing/image/ImageDataGenerator

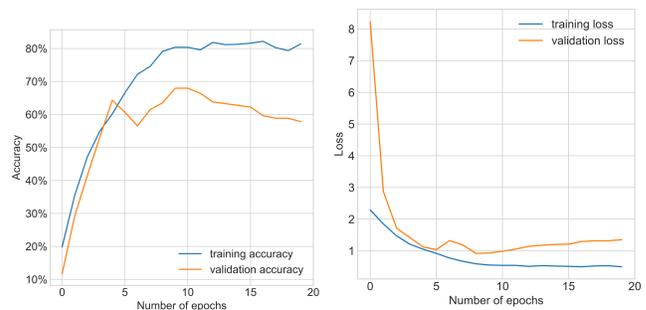


(a) Train dataset.



(b) Test dataset.

Figure 6: Distribution of emotions across the train and test split of the RAF-DB dataset.



(a) Accuracy during training

(b) Loss during training

Figure 7: Accuracy and loss during training the baseline CK+ classifier.

often expressed more subtly, lighting can vary and faces may be obfuscated as illustrated in Figure 5.

Additionally, the RAF-DB dataset provides labels across three demographic categories - **gender** (with subjects labelled as one of “male”, “female” or “unsure”), **race** (“Caucasian”, “African-American” or “Asian”) and **age** (with subjects being assigned to one of 5 age groups - 0-3, 4-19, 20-39, 40-69 and 70+). We form the train and test dataset preserving the original split defined by the authors [31]. This gives us 12271 training images and 3068 test images, with emotions distributed as shown in Figure 6.

3.2 Baseline Models

We use the CK+ and RAF-DB datasets to implement two FER classifiers to serve as baselines to which we will apply the compression strategies. We follow an FER tutorial by S. Kekre [25] to set up a DNN classifier in Keras [14]. The architecture for the CK+ baseline

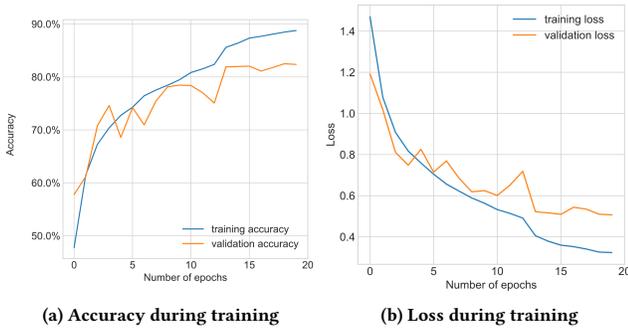


Figure 8: Accuracy and loss during training the baseline RAF-DB classifier.

is shown in Table 1, and the RAF-DB architecture can be found in Table 2 - the two are almost identical except for minor differences to account for the different input size and the different number of classes. The neural architecture is inspired by a study by Goodfellow et al. [17] and contains 4 convolutional layers, followed by 2 hidden fully-connected layers (with pooling and dropout layers in-between). At the time of its publication, the architecture achieved a then state-of-the-art performance of around 65% on the FER-2013 dataset [4, 17].

We compile the baseline models using an Adam optimiser [26] and a categorical cross-entropy loss [18]. We train each model for 20 iterations keeping track of training and validation accuracy, and training and validation loss (where training is performed only on the train dataset and validation is performed only on the test dataset). At every iteration, we store the “best” weights observed so far (i.e., the ones associated with the highest validation accuracy).

The process of obtaining the CK+ and the RAF-DB baseline is illustrated in Figure 7 and Figure 8 respectively. The optimal weights are obtained after the 10th iteration for CK+ and after the 19th iteration for RAF-DB when the two models report respectively 67.96% and 82.46% validation accuracy.

3.3 Model Compression Implementation

We implement three model compression strategies – magnitude-based weight pruning, post-training quantisation and weight clustering. To do this, we make use of TensorFlow’s *Model Optimization Toolkit*⁵, part of the TFLite framework [28].

3.3.1 Quantisation. To quantise the model, we convert the pre-trained Keras baseline described in the last section to a TFLite model and apply the default TFLite optimisation strategy⁶ which reduces the model representation to 8 bits. Finally, we store the quantised model on disk and compress it using the zip compression tool, so that we are able to observe the change in size that quantisation has introduced.

Table 1: Summary of the architecture of the CK+ baseline model generated using Keras’s `model.summary()`. The None values indicate the batch size is flexible.

Layer	Output Shape
conv2d	(None, 48, 48, 64)
batch_normalization	(None, 48, 48, 64)
activation	(None, 48, 48, 64)
max_pooling2d	(None, 24, 24, 64)
dropout	(None, 24, 24, 64)
conv2d_1	(None, 24, 24, 128)
batch_normalization_1	(None, 24, 24, 128)
activation_1	(None, 24, 24, 128)
max_pooling2d_1	(None, 12, 12, 128)
dropout_1	(None, 12, 12, 128)
conv2d_2	(None, 12, 12, 512)
batch_normalization_2	(None, 12, 12, 512)
activation_2	(None, 12, 12, 512)
max_pooling2d_2	(None, 6, 6, 512)
dropout_2	(None, 6, 6, 512)
conv2d_3	(None, 6, 6, 512)
batch_normalization_3	(None, 6, 6, 512)
activation_3	(None, 6, 6, 512)
max_pooling2d_3	(None, 3, 3, 512)
dropout_3	(None, 3, 3, 512)
flatten	(None, 4608)
dense	(None, 256)
batch_normalization_4	(None, 256)
activation_4	(None, 256)
dropout_4	(None, 256)
dense_1	(None, 512)
batch_normalization_5	(None, 512)
activation_5	(None, 512)
dropout_5	(None, 512)
dense_2	(None, 8)
Total trainable parameters:	4,479,240

3.3.2 Pruning. We apply pruning using TFLite’s `ConstantSparsity`⁷ pruning schedule. Our implementation is parameterised by the pruning sparsity – we observe the effect of this parameter on compression in Section 4. After pruning has been applied, we fine-tune the pruned model for 2 iterations as suggested by the TFLite documentation. Similarly to quantisation, we store the model on disk and compress it to evaluate the reduction in size.

3.3.3 Weight Clustering. We implement weight clustering using TFLite’s `cluster_weights`⁸ module and parameterise it by the number of clusters. Similar to pruning, we fine-tune the clustered model, store it on disk and compress it with zip.

3.3.4 Combined Compression. Additionally, our implementation allows applying quantisation on top of a pruned or a weight clustered model. This gives us two additional “hybrid” compression strategies - *pruning with quantisation*, and *weight clustering with*

⁵https://www.tensorflow.org/model_optimization/guide

⁶https://www.tensorflow.org/api_docs/python/tf/lite/Optimize

⁷https://www.tensorflow.org/model_optimization/api_docs/python/tfmot/sparsity/keras/ConstantSparsity

⁸https://www.tensorflow.org/model_optimization/api_docs/python/tfmot/clustering/keras/cluster_weights

quantisation. We use those in our evaluation to explore the combined effect of compression techniques.

4 EXPERIMENTS

In this section we introduce the metrics we have considered during our experiments, and present the results we have obtained. We focus on the more interesting results from the study, however a more detailed breakdown of results is provided in the appendix as Table 6 and Table 7.

4.1 Metrics

In our experiments, we compare the uncompressed baseline model against compressed versions of it across 3 metrics:

- **Model size** – that is the size of the model on disk in megabytes. This metric measures how effective a compression strategy is in reducing the storage requirement of the model. While some of the compression strategies could also reduce other system metrics such as latency, we focus on model size since all three compression techniques are primarily used to reduce storage consumption.

Table 2: Summary of the architecture of the RAF-DB baseline model generated using Keras’s `model.summary()`. The None values indicate the batch size is flexible.

Layer	Output Shape
conv2d	(None, 100, 100, 64)
batch_normalization	(None, 100, 100, 64)
activation	(None, 100, 100, 64)
max_pooling2d	(None, 50, 50, 64)
dropout	(None, 50, 50, 64)
conv2d_1	(None, 50, 50, 128)
batch_normalization_1	(None, 50, 50, 128)
activation_1	(None, 50, 50, 128)
max_pooling2d_1	(None, 25, 25, 128)
dropout_1	(None, 25, 25, 128)
conv2d_2	(None, 25, 25, 512)
batch_normalization_2	(None, 25, 25, 512)
activation_2	(None, 25, 25, 512)
max_pooling2d_2	(None, 12, 12, 512)
dropout_2	(None, 12, 12, 512)
conv2d_3	(None, 12, 12, 512)
batch_normalization_3	(None, 12, 12, 512)
activation_3	(None, 12, 12, 512)
max_pooling2d_3	(None, 6, 6, 512)
dropout_3	(None, 6, 6, 512)
flatten	(None, 18432)
dense	(None, 256)
batch_normalization_4	(None, 256)
activation_4	(None, 256)
dropout_4	(None, 256)
dense_1	(None, 512)
batch_normalization_5	(None, 512)
activation_5	(None, 512)
dropout_5	(None, 512)
dense_2	(None, 7)
Total trainable parameters:	8,013,703

- **Accuracy** – this is the overall accuracy a model achieves on the test dataset. We use this metric as an indicator of how compression has impacted predictive accuracy.
- **Gender accuracy** – to examine the fairness of the models, we also introduce the measure of gender accuracy, specifically as female vs. male accuracy. We define female accuracy as the number of correctly classified images containing a female subject over the total number of images containing a female subject. Similarly, male accuracy is the number of correctly classified images where the subject was male over the total number of images where the subject was male. Under the *overall accuracy equality* formulation of fairness, which we defined in section 2.1, an unbiased model should have equal or similar female and male accuracy metrics. Conversely, a large discrepancy between the model’s accuracy for males and females would be a strong indicator of algorithmic bias. Just like overall accuracy, male and female accuracy are measured on the test dataset.
- **Race accuracy** – this is similar to the gender accuracy, but this time we aim to examine the fairness of the models for different race groups as defined in the corresponding dataset. An unbiased model should have equal or similar accuracy for different race groups. Conversely, a large discrepancy between the model’s accuracy for different race groups would be a strong indicator of algorithmic bias. These are measured on the test dataset.
- **Age accuracy** – this is similar to the gender and ethnicity accuracy, but this time we aim to examine the fairness of the models for different age groups as defined in the corresponding dataset. An unbiased model should have equal or similar accuracy for different age groups. These again are measured on the test dataset.

4.2 CK+ Experiments

4.2.1 Baseline Performance. We first report the baseline performance against which we compare the performance of the compressed models. As mentioned previously, the baseline model reports an overall accuracy of **67.96%** on the test dataset. We measure the baseline’s size in the same way we measure the size of the compressed models – we save the model as a file on disk and apply zip compression to the file. After doing that, we find that the baseline’s model size on disk is 16,512,048 bytes or around **16.51 megabytes**.

We find that the female accuracy of the baseline model is **67.08%** and the male accuracy is **69.44%**. This is interesting because as we mentioned in Section 3.1, the CK+ dataset is imbalanced in favour of female subjects and therefore we would expect the baseline model to classify females more accurately.

One reason why the classifier might perform slightly better on male faces is the slight difference in the distribution of emotions across males and females in CK+ – for instance, only 2.1% of male subjects have expressed “contempt” while female subjects have expressed this emotion more than twice more frequently (5.02%). If contempt is an emotion that is inherently harder to classify, then this difference could translate into a minor advantage for classifying male subjects. In any case, though, the gap between male and female accuracy is too minor to conclude the baseline model is biased.

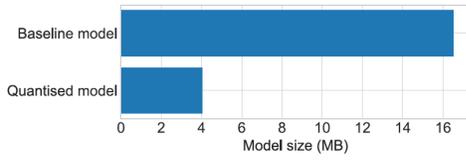


Figure 9: Size of the CK+ classifier before (top) and after (bottom) applying quantisation.

4.2.2 Quantisation Results. We evaluate quantisation by quantising the baseline model 3 times and reporting the mean values for each metric.

After applying quantisation to the baseline model, we observe a 4× reduction in the model size as illustrated in Figure 9. Moreover, this compression comes at no cost – there is no change in the predictive accuracy or fairness whatsoever: Overall accuracy, female accuracy and male accuracy have all remained completely identical to those of the baseline model. Quantisation therefore preserves the fairness and predictive accuracy of the model, while introducing a significant reduction in its size.

4.2.3 Pruning Results. We evaluate the pruning strategy at 6 different levels of sparsity: 10%, 20%, 30%, 40%, 50% and 60%. For each level of sparsity, we prune the baseline model three times and report the mean values for each metric.

We can observe the trade-off that pruning sparsity introduces in Figure 10: As sparsity increases, the model size decreases linearly, reducing the model size by a half at 60% sparsity. However, high sparsity also reduces the network capacity which can impact the accuracy of the model [37]. We can see in Figure 10 (b) that overall accuracy steadily starts to decrease at 40% and 50% sparsity before plummeting at 60%.

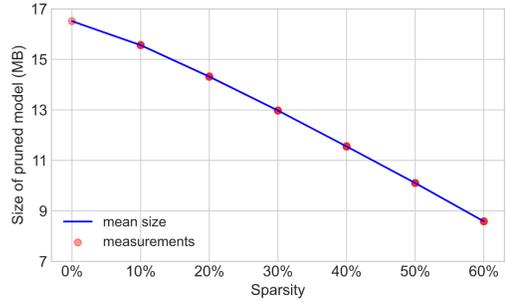
Except for the considerable drop of accuracy at 60% sparsity, though, we can conclude that pruning has had little impact on overall accuracy: At 50% pruning, accuracy has only dropped from 67.96% down to 64.06%.

However, despite the minor drop in overall accuracy, we find that pruning has dramatically increased the discrepancy between female and male accuracy: Table 3 shows that pruning the model tends to keep male accuracy high (in fact, male accuracy has increased for all sparsities up to 50%) while deteriorating female accuracy (which has dropped by 7.22% at 50% sparsity).⁹

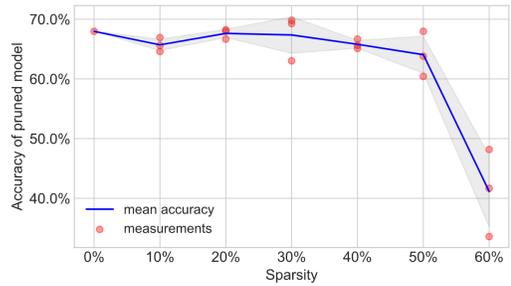
Those results are in agreement with the study by Hooker et al., which finds that “minimal changes to overall accuracy hide disproportionately high errors” [22] in subgroups. Pruning the baseline model by up to 50% has led to a considerable reduction in size with seemingly low impact on overall accuracy. However, while overall accuracy has stayed reasonably high, the initial 2.36% gap between male and female accuracy has grown to 15.6%, amplifying the minimal bias in the baseline model.

As mentioned previously, we are also interested in the “combined” effect of compression techniques. To this end, we apply quantisation on top of the pruned models. We find that quantisation

⁹Interestingly, pruning at 60% actually reports better female than male accuracy. However, at 60% sparsity both male and female accuracy are too low for such a model to be of practical interest.



(a) Model size after pruning.



(b) Model accuracy after pruning. Standard deviation is shaded in gray.

Figure 10: Pruning’s effect on model size and accuracy for CK+ DB.

Table 3: Pruning’s effect on fairness for CK+ DB.

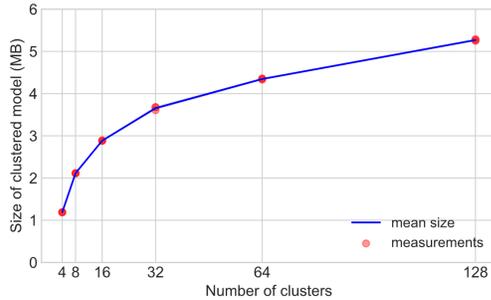
Sparsity	Female accuracy	Male accuracy	Gap
0% (baseline)	67.08%	69.44%	2.36%
10%	59.86%	75.46%	15.60%
20%	62.77%	75.69%	12.92%
30%	62.36%	75.69%	13.33%
40%	63.47%	69.67%	6.20%
50%	59.86%	71.06%	11.20%
60%	44.02%	36.34%	7.68%

can greatly enhance the compression effect of pruning: Quantising the pruned models has decreased their size by a further 3.5 times on average (exact results reported in Table 6 in the appendix).

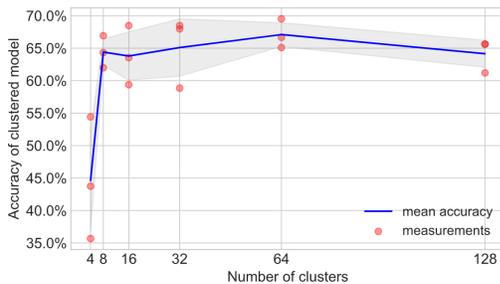
Meanwhile, quantisation has not changed the overall accuracy, male accuracy or female accuracy of the pruned models by more than 1%. Again, applying quantisation is “for free” since no significant impact on predictive performance or fairness is observed.

4.2.4 Weight Clustering Results. We evaluate weight clustering with 4, 8, 16, 32, 64 and 128 clusters. For each number of clusters, we run weight clustering 3 times and report the mean values for each metric.

Similar to pruning, the “number of clusters” parameter introduces a size-accuracy trade-off illustrated in Figure 11. Decreasing the number of clusters rapidly decreases the baseline model size



(a) Model size after clustering.



(b) Model accuracy after weight clustering. Standard deviation is shaded in gray.

Figure 11: Weight clustering’s effect on size and overall accuracy for CK+ DB.

(shrinking it by almost 14 times at 4 clusters). However, an excessively low number of clusters can decrease accuracy dramatically (with overall accuracy dropping below 45% at 4 clusters).

Despite that, we can see that keeping the number of clusters sufficiently high (between 8 and 128) preserves overall accuracy close to the baseline accuracy of 67.96%, while offering a lucrative reduction in model size.

Just like with pruning, though, the overall accuracy of the clustered models is deceptive as it hides an increasing gap between male and female accuracy. Table 4 shows that the baseline 2.36% gap has increased massively, reaching 19.77% in favour of male subjects at 16 clusters.

Applying quantisation to the clustered models reduces their size by a further 17% on average – a much smaller reduction than was observed for pruning. Again, though, quantisation comes at no cost for accuracy or fairness – overall accuracy, male accuracy and female accuracy have all stayed within 1.5% of the values for the clustered models.

4.3 RAF-DB Experiments

We now present the results obtained on the RAF-DB dataset. We provide the full results in Table 7 in the appendix, and only summarise the more interesting results below.

4.3.1 Baseline Performance. The performance of the baseline RAF-DB classifier is summarised in Table 5. The baseline model has a size of **29.80 MB** and reports overall test accuracy of **82.46%**, which

Table 4: Weight clustering’s effect on fairness for CK+ DB.

Number of clusters	Female accuracy	Male accuracy	Gap
4	43.05%	47.22%	4.17%
8	59.30%	72.91%	13.61%
16	56.38%	76.15%	19.77%
32	60.55%	72.68%	12.13%
64	60.27%	78.47%	18.20%
128	59.16%	72.45%	13.29%

Table 5: Metrics for the baseline RAF-DB classifier.

Metric	Value
Size	29.80 MB
Overall accuracy	82.46%
Female accuracy	83.33%
Male accuracy	80.54%
Caucasian accuracy	81.92%
African-American accuracy	86.75%
Asian accuracy	83.02%
A0 accuracy	89.96%
A1 accuracy	82.96%
A2 accuracy	80.44%
A3 accuracy	85.85%
A4 accuracy	70.78%

seems acceptable given that much larger, state-of-the-art architectures have reported between 86% and 89% on this dataset [56]. In terms of fairness, the classifier seems to classify female subjects slightly more accurately than male subjects, and African-American subjects better than Caucasian or Asian subjects. However, those differences are minor. A more major classification gap is observed with respect to the age attribute where there is a 19% classification gap between the best classified age group (A0 or 0-3 years old) and the worst classified age group (A4 or 70+ years old). This could be due to younger subjects expressing emotions more explicitly, or a different distribution of emotion labels across the two age groups. Such a conclusion is supported by relevant works indicating that the age of the face plays an important role for facial expression decoding and factors such as lower expressivity and age-related changes in the face may lower decoding accuracy for older faces [16].

4.3.2 Quantisation Results. Figure 12 shows the size of the RAF-DB classifier before (top) and after (bottom) applying quantisation. We observe that applying quantisation to the baseline classifier reduces the size of the model by around 4.5 times - from 29.80 MB down to 6.56 MB. At the same time, similar to the CK+ experiments, applying quantisation does not impact the predictive performance of the classifier. The quantised model’s overall accuracy is identical to the one reported by the baseline model (82.46%), and Table 7 shows that none of the per-class accuracies have changed by more than 0.2%.

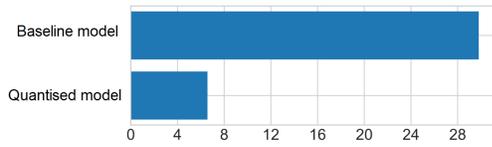


Figure 12: Size of the RAF-DB classifier before (top) and after (bottom) applying quantisation.

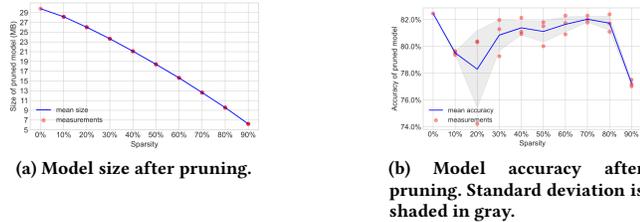


Figure 13: Pruning’s effect on model size and accuracy for RAF-DB.

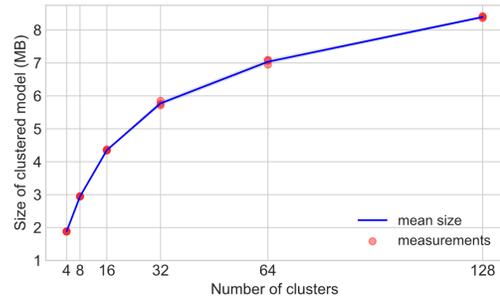
4.3.3 Pruning Results. Similar to the CK+ experiments, we observe that applying pruning to the RAF-DB classifier translates to a linear reduction in model size as illustrated in Figure 13. Unlike the CK+ classifier though, the accuracy of the RAF-DB model seems to be much more robust to pruning. Even at 80% sparsity, the overall accuracy has only dropped down to 81.73% compared to the baseline (82.46%). Only when sparsity increases to 90% do we observe a more significant drop in accuracy down to 77.23%. The RAF-DB model is therefore more akin to models such as MNIST classifiers [5] where near-optimal accuracy can be preserved even at 99% sparsity. We analyse the difference in robustness between the CK+ and the RAF-DB model, and discuss potential causes for this disparity in the Discussion section.

As for gender fairness, the original 2.8% classification gap varies between 0.2% and 2.7% depending on the level of sparsity, so there is no evidence suggesting that sparsity has negatively impacted fairness. Since RAF-DB also has race and age labels, we can analyse fairness across those dimensions as well. The full fairness metrics are presented in Table 7, which shows that applying pruning to the RAF-DB classifier has little to no effect on model age and race fairness too.

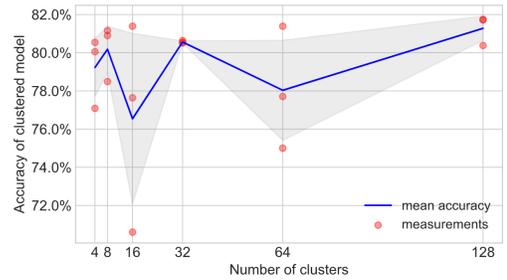
4.3.4 Weight Clustering Results. Figure 13 shows pruning’s effect on model size and accuracy for RAF-DB. In terms of model accuracy, the classifier trained and tested on RAF-DB seems more robust to compression compared to the CK+ experiments. Regardless of the level of sparsity or the number of clusters of the compression strategy, the overall accuracy of the compressed model does not fall below 77%, which is not significantly lower than the uncompressed (or baseline) model, which has an accuracy of 82%.

5 DISCUSSION

In this study, we analysed and compared the effect of model compression on model size, accuracy and fairness in the context of



(a) Model size after clustering.



(b) Model accuracy after weight clustering. Standard deviation is shaded in gray.

Figure 14: Weight clustering’s effect on size and overall accuracy for RAF DB.

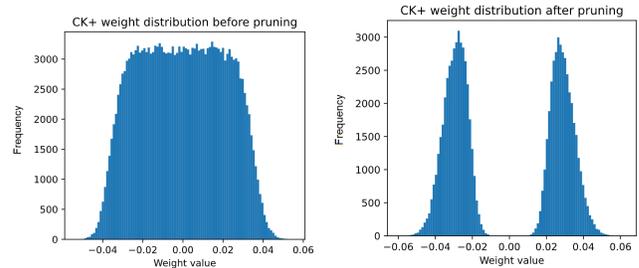


Figure 15: Distribution of the kernel weights of the conv2d layer of the CK+ classifier before (left) and after (right) pruning at 60%.

facial expression recognition on two facial expression datasets. We now revisit our research questions from Section ??:

- **RQ1: “How effective is model compression in the context of FER?”** We saw that model compression can dramatically reduce the storage requirements of both FER models. Quantisation alone achieves around 4 – 4.5× reduction in model sizes with minimal impact on overall accuracy for both CK+ and RAF-DB. Both pruning and clustering progressively decrease the model size on disk as sparsity increases (and the number of clusters decrease). In terms of model accuracy, the classifier trained and tested on RAF-DB seems more robust

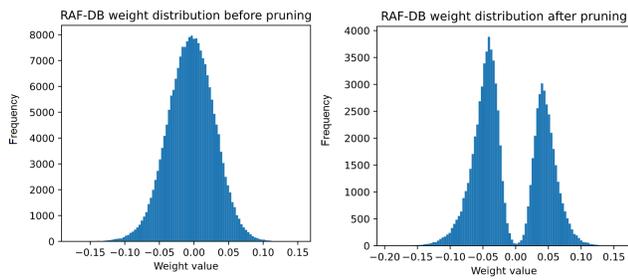


Figure 16: Distribution of the kernel weights of the conv2d layer of the RAF-DB classifier before (left) and after (right) pruning at 60%.

to compression compared to the CK+ DB one. Regardless of the level of sparsity or the number of clusters of the compression strategy, the overall accuracy of the compressed model does not fall below 77%.

- **RQ2: “Do model compression techniques amplify biases?”** Our findings for CK+ DB confirm the claims by previous studies that model compression can amplify existing biases for gender for deep learning models trained for facial expression recognition. However our findings for RAF-DB indicate that sparsity does not impact fairness in terms of gender, race or age negatively.
- **RQ3: “Is the impact on fairness identical across different compression techniques?”** Our results for CK+ DB suggest that different compression techniques tend to have a highly distinct impact on fairness: Post-training quantisation has no visible effect on baseline fairness. Meanwhile, pruning and weight clustering can severely amplify biases, increasing the initial 2.36% gap between male and female accuracy to 15.60% and 19.77% respectively. On the other hand for RAF-DB we find that the different compression strategies do not seem to increase the gap in predictive performance across any of the three demographic attributes (gender, race and age). That is in contrast with the CK+ experiment findings where compression seems to amplify existing biases.

In order to understand the reasons for the different findings related to these datasets, we compare the weight distributions of CK+ and RAF-DB before pruning, and we see that the CK+ distribution has a ‘wider’ shape compared to the RAF-DB one which is much more ‘narrow’ and most of its values are located close to its mean at 0.00. As a result, when we prune the two models, we get gaps of different sizes. For the ‘wide’ CK+ distribution, when we prune at 60%, we need to set ‘crop’ or zero-out 60% of its weights. However, the CK+ weights are relatively evenly distributed and most of them are located at some (relative) distance from the 0.00 mean. As a result, we need to ‘crop’ weights which are not located immediately around the centre and that causes a wide gap in the middle. We can expect that a large gap will have a bigger impact on predictive performance since it means values which are further from 0.00 have been set to 0.00. Meanwhile, for the RAF-DB the resulting gap is much smaller. This is because the original distribution is much

more ‘narrow’ in terms of shape with a big share of the weights located at or close to 0.00. Therefore, when we prune at 60%, the 60% of the weights which we set to 0.00 are going to be already equal to or close to 0.00 and therefore we can expect a minor impact on predictive performance. This can explain why RAF-DB is much more robust than CK+ to pruning. As a reminder, at 60% pruning, CK+’s accuracy dropped from 68% down to 42%, and RAF-DB’s accuracy only dropped from 82.4% down to 81.6%. While the plots in this paper only show the weights from the first dense layer, we observe a similar trend for the weights of the other dense layers, as well as the weights of the convolutional layers. It is important to note that the CK+DB is much more homogenous in terms of acquisition setup and setting as compared to RAF-DB which contains various facial images crawled from the internet, and the size of CK+DB is significantly smaller than RAF-DB.

6 LIMITATIONS AND FUTURE WORK

We identify a couple of **limitations** of our study: First, the gender annotation of the CK+ dataset was performed manually. While the annotations were straightforward, labels could be obtained via a crowdsourcing experiment or a user study to better ensure the ground truth is reliable and unbiased.

Furthermore, our baseline model is not a state-of-the-art FER classifier as of 2021. We deliberately selected a relatively small architecture that could be trained, fine-tuned and evaluated locally, but given more time and computational resources the study could be extended to consider a larger and more modern architecture.

The work in this paper could be extended in several different directions: More compression strategies (e.g., quantisation-aware training and various forms of weight sharing [29]) could also be evaluated. The study could also be extended to explore more system metrics such as latency, memory consumption, etc.

7 ACKNOWLEDGMENTS

S. Stoychev completed this work while studying towards his MPhil in Advanced Computing in the Department of Computer Science and Technology, University of Cambridge. H. Gunes is supported by the EPSRC project ARoEQ under grant ref. EP/R030782/1.

REFERENCES

- [1] 2019. Apple’s ‘sexist’ credit card investigated by US regulator. <https://www.bbc.co.uk/news/business-50365609>
- [2] 2019. Introducing AI Fairness 360, A Step Towards Trusted AI - IBM Research. <https://www.ibm.com/blogs/research/2018/09/ai-fairness-360/>
- [3] 2020. TensorFlow Model Optimization Toolkit - Weight Clustering API. <https://blog.tensorflow.org/2020/08/tensorflow-model-optimization-toolkit-weight-clustering-api.html>
- [4] Abhinav Agrawal and Namita Mittal. 2020. Using CNN for facial expression recognition: a study of the effects of kernel size and number of filters on accuracy. *The Visual Computer* 36 (02 2020). <https://doi.org/10.1007/s00371-019-01630-9>
- [5] Simon Alford, Ryan Robinett, Lauren Milechin, and Jeremy Kepner. 2019. Training behavior of sparse neural network topologies. In *2019 IEEE High Performance Extreme Computing Conference (HPEC)*. IEEE, 1–6.
- [6] Pinar Barlas, Kyriakos Kyriakou, Styliani Kleantous, and Jahna Otterbacher. 2019. Social B(eye)s: Human and Machine Descriptions of People Images. *Proceedings of the International AAI Conference on Web and Social Media* 13, 01 (Jul. 2019), 583–591. <https://ojs.aaai.org/index.php/ICWSM/article/view/3255>
- [7] Solon Barocas, Moritz Hardt, and Arvind Narayanan. 2019. *Fairness and Machine Learning*. fairmlbook.org. <http://www.fairmlbook.org>
- [8] M.S. Bartlett, G. Littlewort, M. Frank, C. Lainscsek, I. Fasel, and J. Movellan. 2005. Recognizing facial expression: machine learning and application to spontaneous behavior. In *2005 IEEE Computer Society Conference on Computer Vision and*

- Pattern Recognition (CVPR'05)*, Vol. 2. 568–573 vol. 2. <https://doi.org/10.1109/CVPR.2005.297>
- [9] Davis Blalock, Jose Javier Gonzalez Ortiz, Jonathan Frankle, and John Guttag. 2020. What is the State of Neural Network Pruning? arXiv:2003.03033 [cs.LG].
 - [10] Joy Buolamwini and Timnit Gebru. 2018. Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification. In *FAT*.
 - [11] Maurizio Capra, Beatrice Bussolino, Alberto Marchisio, Muhammad Shafique, Guido Masera, and Maurizio Martina. 2020. An Updated Survey of Efficient Hardware Architectures for Accelerating Deep Convolutional Neural Networks. *Future Internet* 12, 7 (2020). <https://doi.org/10.3390/fi12070113>
 - [12] Yu Cheng, Duo Wang, Pan Zhou, and Tao Zhang. 2017. A Survey of Model Compression and Acceleration for Deep Neural Networks. *CoRR abs/1710.09282* (2017). arXiv:1710.09282 <http://arxiv.org/abs/1710.09282>
 - [13] Jiace Cheong, Sinan Kalkan, and Hatice Gunes. 2021. The Hitchhiker’s Guide to Bias and Fairness in Facial Affective Signal Processing: Overview and techniques. *IEEE Signal Process. Mag.* 38, 6 (2021), 39–49. <https://doi.org/10.1109/MSP.2021.3106619>
 - [14] François Chollet. 2015. Keras. <https://github.com/fchollet/keras>.
 - [15] Jeffrey Dastin. 2018. Amazon scraps secret AI recruiting tool that showed bias against women. <https://www.reuters.com/article/us-amazon-com-jobs-automation-insight-idUSKCN1MK08G>
 - [16] Mara Folster, Ursula Hess, and Katja Werheid. 2014. Facial age affects emotional expression decoding. *Frontiers in Psychology* 5:30 (2014), 1–13.
 - [17] Ian J. Goodfellow, Dumitru Erhan, Pierre Luc Carrier, Aaron Courville, Mehdi Mirza, Ben Hamner, Will Cukierski, Yichuan Tang, David Thaler, Dong-Hyun Lee, Yingbo Zhou, Chetan Ramaiah, Fangxiang Feng, Ruifan Li, Xiaojie Wang, Dimitris Athanasakis, John Shawe-Taylor, Maxim Milakov, John Park, Radu Ionescu, Marius Popescu, Cristian Groza, James Bergstra, Jingjing Xie, Lukasz Romaszko, Bing Xu, Zhang Chuang, and Yoshua Bengio. 2013. Challenges in Representation Learning: A report on three machine learning contests. arXiv:1307.0414 [stat.ML]
 - [18] Elliott Gordon-Rodriguez, Gabriel Loaiza-Ganem, Geoff Pleiss, and John P. Cunningham. 2020. Uses and Abuses of the Cross-Entropy Loss: Case Studies in Modern Deep Learning. arXiv:2011.05231 [stat.ML]
 - [19] Yuanyuan Guo, Yifan Xia, Jing Wang, Hui Yu, and Rung-Ching Chen. 2020. Real-time facial affective computing on mobile devices. *Sensors* 20, 3 (6 Feb. 2020). <https://doi.org/10.3390/s20030870>
 - [20] Song Han, Huizi Mao, and William J. Dally. 2015. Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding. (2015). <http://arxiv.org/abs/1510.00149> cite arxiv:1510.00149Comment: Published as a conference paper at ICLR 2016 (oral).
 - [21] Will Douglas Heaven. 2020. Predictive policing algorithms are racist. They need to be dismantled. <https://www.technologyreview.com/2020/07/17/1005396/predictive-policing-algorithms-racist-dismantled-machine-learning-bias-criminal-justice/>
 - [22] Sara Hooker, Nyalleng Moorosi, Gregory Clark, Samy Bengio, and Emily Denton. 2020. Characterising Bias in Compressed Models. arXiv:2010.03058 [cs.LG]
 - [23] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew G. Howard, Hartwig Adam, and Dmitry Kalenichenko. 2017. Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference. *CoRR abs/1712.05877* (2017). arXiv:1712.05877 <http://arxiv.org/abs/1712.05877>
 - [24] Khari Johnson. 2020. Algorithmic Justice League protests bias in voice AI and media coverage. <https://venturebeat.com/2020/03/31/algorithmic-justice-league-protests-bias-voice-ai-and-media-coverage/>
 - [25] Snehan Kekre. [n.d.]. Tutorial: Facial Expression Recognition with Keras. <https://www.coursera.org/projects/facial-expression-recognition-keras>.
 - [26] Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). <http://arxiv.org/abs/1412.6980>
 - [27] Andreas Kleppe, Ole-Johan Skrede, Sepp De Raedt, Knut Liestøl, David J Kerr, and Håvard E Danielsen. 2021. Designing deep learning studies in cancer diagnostics. *Nature reviews. Cancer* 21, 3 (March 2021), 199–211. <https://doi.org/10.1038/s41568-020-00327-9>
 - [28] Juhyun Lee, Nikolay Chirkov, Ekaterina Ignasheva, Yury Pisarchyk, Mogan Shieh, Fabio Riccardi, Raman Sarokin, Andrei Kulik, and Matthias Grundmann. 2019. On-Device Neural Net Inference with Mobile GPUs. *CoRR abs/1907.01989* (2019). arXiv:1907.01989 <http://arxiv.org/abs/1907.01989>
 - [29] Seulki Lee and Shahriar Nirjon. 2020. Fast and Scalable In-Memory Deep Multi-task Learning via Neural Weight Virtualization. In *Proceedings of the 18th International Conference on Mobile Systems, Applications, and Services* (Toronto, Ontario, Canada) (*MobiSys '20*). Association for Computing Machinery, New York, NY, USA, 175–190. <https://doi.org/10.1145/3386901.3388947>
 - [30] Jing-Ping Li, Nawazish Mirza, Birjees Rahat, and Deping Xiong. 2020. Machine learning and credit ratings prediction in the age of fourth industrial revolution. *Technological Forecasting and Social Change* 161 (2020), 120309. <https://doi.org/10.1016/j.techfore.2020.120309>
 - [31] Shan Li, Weihong Deng, and JunPing Du. 2017. Reliable Crowdsourcing and Deep Locality-Preserving Learning for Expression Recognition in the Wild. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2584–2593.
 - [32] Tailin Liang, John Glossner, Lei Wang, and Shaobo Shi. 2021. Pruning and Quantization for Deep Neural Network Acceleration: A Survey. arXiv:2101.09671 [cs.CV]
 - [33] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. 2015. Deep Learning Face Attributes in the Wild. In *Proceedings of International Conference on Computer Vision (ICCV)*.
 - [34] Patrick Lucey, Jeffrey F. Cohn, Takeo Kanade, Jason Saragih, Zara Ambadar, and Iain Matthews. [n.d.]. The Extended Cohn-Kanade Dataset (CK+): A complete dataset for action unit and emotion-specified expression.
 - [35] Yury Nahshan, Brian Chmiel, Chaim Baskin, Evgenii Zheltonozhskii, Ron Banner, Alexander M. Bronstein, and Avi Mendelson. 2019. Loss Aware Post-training Quantization. *CoRR abs/1911.07190* (2019). arXiv:1911.07190 <http://arxiv.org/abs/1911.07190>
 - [36] Luca Oneto and Silvia Chiappa. 2020. *Fairness in Machine Learning*. Springer International Publishing, Cham, 155–196. https://doi.org/10.1007/978-3-030-43883-8_7
 - [37] Michela Paganini. 2020. Prune Responsibly. arXiv:2009.09936 [cs.CV]
 - [38] Omkar M. Parkhi, Andrea Vedaldi, and Andrew Zisserman. 2015. Deep Face Recognition. In *Proceedings of the British Machine Vision Conference (BMVC)*, Xianghua Xie, Mark W. Jones, and Gary K. L. Tam (Eds.). BMVA Press, Article 41, 12 pages. <https://doi.org/10.5244/C.29.41>
 - [39] Dana Pessach and Erez Shmueli. 2020. Algorithmic Fairness. arXiv:2001.09784 [cs.CY]
 - [40] Morgan Klaus Scheuerman, Jacob M. Paul, and Jed R. Brubaker. 2019. How Computers See Gender: An Evaluation of Gender Classification in Commercial Facial Analysis Services. *Proc. ACM Hum.-Comput. Interact.* 3, CSCW, Article 144 (Nov. 2019), 33 pages. <https://doi.org/10.1145/3359246>
 - [41] Roy Schwartz, Jesse Dodge, Noah A. Smith, and Oren Etzioni. 2019. Green AI. *CoRR abs/1907.10597* (2019). arXiv:1907.10597 <http://arxiv.org/abs/1907.10597>
 - [42] Jin Shuting, Xiangxiang Zeng, Feng Xia, Wei Huang, and Xiangrong Liu. 2020. Application of deep learning methods in biological networks. *Briefings in bioinformatics* 22 (05 2020). <https://doi.org/10.1093/bib/bbaa043>
 - [43] Matteo Spezialetti, Giuseppe Placidi, and Silvia Rossi. 2020. Emotion Recognition for Human-Robot Interaction: Recent Advances and Future Perspectives. *Frontiers in Robotics and AI* 7 (2020), 145. <https://doi.org/10.3389/frobt.2020.532279>
 - [44] Nate Swanner. 2018. Google’s New Machine Learning Curriculum Aims to Stop Bias Cold. <https://insights.dice.com/2018/10/24/google-machine-learning-course-bias/>
 - [45] Tomasz Szydio, Joanna Sendorek, and Robert Brzoza-Woch. 2018. Enabling Machine Learning on Resource Constrained Devices by Source Code Generation of the Learned Models. In *Computational Science – ICCS 2018*, Yong Shi, Haohuan Fu, Yingjie Tian, Valeria V. Krzhizhanovskaya, Michael Harold Lees, Jack Dongarra, and Peter M. A. Sloot (Eds.). Springer International Publishing, Cham, 682–694.
 - [46] Ameet Talwalkar. 2020. AI in the 2020s Must Get Greener—and Here’s How. <https://spectrum.ieee.org/energywise/artificial-intelligence/machine-learning/energy-efficient-green-ai-strategies>
 - [47] TensorFlow. 2019. TensorFlow Model Optimization Toolkit - Pruning API. <https://medium.com/tensorflow/tensorflow-model-optimization-toolkit-pruning-api-42cac9157a6a>
 - [48] Sahil Verma and Julia Rubin. 2018. Fairness Definitions Explained. In *Proceedings of the International Workshop on Software Fairness* (Gothenburg, Sweden) (*FairWare '18*). Association for Computing Machinery, New York, NY, USA, 1–7. <https://doi.org/10.1145/3194770.3194776>
 - [49] Kai Wang, Xiaojiang Peng, Jianfei Yang, Shijian Lu, and Yu Qiao. 2020. Suppressing Uncertainties for Large-Scale Facial Expression Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
 - [50] Tian Xu, Jennifer White, Sinan Kalkan, and Hatice Gunes. 2020. Investigating Bias and Fairness in Facial Expression Recognition. In *Computer Vision – ECCV 2020 Workshops – Glasgow, UK, August 23–28, 2020, Proceedings, Part VI* 506–523.
 - [51] Tian Xu, Jennifer White, Sinan Kalkan, and Hatice Gunes. 2020. Investigating Bias and Fairness in Facial Expression Recognition. In *Computer Vision – ECCV 2020 Workshops*, Adrien Bartoli and Andrea Fusiello (Eds.). Springer International Publishing, Cham, 506–523.
 - [52] Haibin Yan, Marcelo H. Ang, and Aun Neow Poo. 2011. Cross-dataset facial expression recognition. In *2011 IEEE International Conference on Robotics and Automation*. 5985–5990. <https://doi.org/10.1109/ICRA.2011.5979705>
 - [53] Shen Yan, Di Huang, and Mohammad Soleymani. 2020. Mitigating Biases in Multimodal Personality Assessment. In *Proceedings of the 2020 International Conference on Multimodal Interaction* (Virtual Event, Netherlands) (*ICMI '20*). Association for Computing Machinery, New York, NY, USA, 361–369. <https://doi.org/10.1145/3382507.3418889>
 - [54] I-Cheng Yeh and Che hui Lien. 2009. The comparisons of data mining techniques for the predictive accuracy of probability of default of credit card clients. *Expert Systems with Applications* 36, 2, Part 1 (2009), 2473–2480. <https://doi.org/10.1016/j.eswa.2007.12.020>
 - [55] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, and Krishna P. Gummadi. 2017. Fairness Beyond Disparate Treatment & Disparate Impact.

Proceedings of the 26th International Conference on World Wide Web (Apr 2017).
<https://doi.org/10.1145/3038912.3052660>

- [56] Hengshun Zhou, Debin Meng, Yuanyuan Zhang, Xiaojiang Peng, Jun Du, Kai Wang, and Yu Qiao. 2019. Exploring Emotion Features and Fusion Strategies for Audio-Video Emotion Recognition. *2019 International Conference on Multimodal Interaction (Oct 2019)*. <https://doi.org/10.1145/3340555.3355713>
- [57] Michael Zhu and Suyog Gupta. 2017. To prune, or not to prune: exploring the efficacy of pruning for model compression. arXiv:1710.01878 [stat.ML]

A FULL CK+ RESULTS

Table 6: Full results from the CK+ experiments.

Model	Size (MB)	Overall acc.	Female acc.	Male acc.
baseline	16.51	67.96%	67.08%	69.44%
quantised	4.04	67.96%	67.08%	69.44%
pruned (10%)	15.56	65.71%	59.86%	75.46%
pruned (20%)	14.31	67.62%	62.77%	75.69%
pruned (30%)	12.97	67.36%	62.36%	75.69%
pruned (40%)	11.55	65.79%	63.47%	69.67%
pruned (50%)	10.10	64.06%	59.86%	71.06%
pruned (60%)	8.58	41.11%	44.02%	36.34%
pruned (10%) + quant.	3.88	65.45%	59.58%	75.23%
pruned (20%) + quant.	3.73	67.53%	62.50%	75.92%
pruned (30%) + quant.	3.51	67.36%	62.36%	75.69%
pruned (40%) + quant.	3.21	66.05%	63.88%	69.67%
pruned (50%) + quant.	2.85	64.06%	59.86%	71.06%
pruned (60%) + quant.	2.44	40.88%	43.61%	36.34%
clustered (4 cl.)	1.18	44.61%	43.05%	47.22%
clustered (8 cl.)	2.11	64.40%	59.30%	72.91%
clustered (16 cl.)	2.88	63.80%	56.38%	76.15%
clustered (32 cl.)	3.65	65.10%	60.55%	72.68%
clustered (64 cl.)	4.34	67.10%	60.27%	78.47%
clustered (128 cl.)	5.26	64.14%	59.16%	72.45%
clust. (4 cl.) + quant.	0.91	43.57%	42.08%	46.06%
clust. (8 cl.) + quant.	1.81	64.49%	59.72%	72.45%
clust. (16 cl.) + quant.	2.57	64.49%	57.08%	76.85%
clust. (32 cl.) + quant.	3.22	65.71%	61.25%	73.14%
clust. (64 cl.) + quant.	3.63	66.75%	60.00%	78.00%
clust. (128 cl.) + quant.	3.79	64.49%	59.72%	72.45%

B FULL RAF-DB RESULTS

Table 7: Full results from the RAF-DB experiments. “Female” and “Male” indicate accuracies for female and male subjects respectively. “Cauc.,” “Af.-Am.” and “Asian” denote accuracies for subjects labelled as Caucasian, African-American and Asian respectively. “A0” to “A4” indicate accuracies across the 5 age groups - 0-3, 4-19, 20-39, 40-69 and 70+.

Model	Size (MB)	Overall acc.	Female	Male	Cauc.	Af.-Am.	Asian	A0	A1	A2	A3	A4
baseline	29.80	82.46%	83.33%	80.54%	81.92%	86.75%	83.02%	89.96%	82.92%	80.44%	85.85%	70.78%
quantised	6.56	82.46%	83.20%	80.70%	81.92%	86.75%	83.02%	89.96%	82.92%	80.50%	85.65%	70.78%
pruned (10%)	28.15	79.51%	80.04%	78.64%	79.00%	81.90%	80.88%	82.57%	82.64%	78.53%	79.88%	67.41%
pruned (20%)	26.01	78.30%	78.72%	77.47%	78.10%	80.19%	78.32%	82.16%	79.35%	77.47%	79.34%	67.79%
pruned (30%)	23.65	80.84%	81.37%	79.69%	80.44%	84.33%	81.09%	85.51%	82.71%	78.90%	84.32%	70.03%
pruned (40%)	21.09	81.38%	81.95%	79.93%	80.97%	84.04%	82.10%	88.34%	82.78%	79.56%	83.79%	68.53%
pruned (50%)	18.41	81.11%	81.83%	79.47%	80.85%	84.33%	80.81%	87.03%	82.85%	79.18%	83.59%	71.91%
pruned (60%)	15.63	81.64%	82.53%	79.82%	81.18%	83.76%	82.88%	87.03%	82.57%	80.42%	82.73%	73.40%
pruned (70%)	12.62	82.02%	82.40%	80.62%	81.75%	84.90%	81.98%	87.84%	83.88%	80.22%	84.32%	71.16%
pruned (80%)	9.55	81.73%	82.26%	80.11%	81.35%	84.90%	82.05%	88.34%	83.19%	79.88%	83.86%	71.91%
pruned (90%)	6.21	77.23%	77.09%	76.83%	76.77%	78.20%	79.02%	82.26%	78.87%	75.35%	80.34%	67.41%
pruned (10%) + quant.	6.44	79.54%	80.06%	78.62%	79.05%	81.90%	80.74%	82.67%	82.57%	78.58%	79.94%	67.04%
pruned (20%) + quant.	6.29	78.32%	78.68%	77.58%	78.16%	80.19%	78.19%	82.06%	79.49%	77.45%	79.54%	67.41%
pruned (30%) + quant.	6.03	80.85%	81.44%	79.58%	80.50%	84.33%	80.88%	85.71%	82.78%	78.92%	84.19%	69.66%
pruned (40%) + quant.	5.55	81.46%	82.05%	79.98%	81.05%	84.04%	82.19%	88.34%	82.78%	79.68%	83.86%	68.53%
pruned (50%) + quant.	5.02	81.04%	81.83%	79.47%	80.78%	84.33%	80.67%	87.03%	82.71%	79.18%	83.20%	72.28%
pruned (60%) + quant.	4.38	81.64%	82.46%	79.90%	81.19%	83.76%	82.81%	86.93%	82.51%	80.42%	82.86%	73.40%
pruned (70%) + quant.	3.57	82.05%	82.42%	80.65%	81.79%	84.75%	81.98%	87.84%	84.01%	80.26%	84.19%	71.16%
pruned (80%) + quant.	2.73	81.70%	82.26%	80.11%	81.36%	84.75%	81.84%	88.14%	83.19%	79.76%	84.19%	71.91%
pruned (90%) + quant.	1.65	77.33%	77.26%	76.88%	76.88%	77.92%	79.22%	82.16%	78.94%	75.49%	80.47%	67.41%
clustered (4 cl.)	1.88	79.22%	79.71%	78.08%	79.05%	79.05%	80.12%	85.61%	81.55%	76.93%	81.80%	71.16%
clustered (8 cl.)	2.94	80.18%	80.39%	79.34%	79.75%	81.48%	81.64%	85.10%	81.61%	78.68%	82.27%	70.41%
clustered (16 cl.)	4.35	76.54%	77.44%	75.02%	75.92%	78.77%	78.46%	80.64%	79.21%	75.15%	78.55%	61.42%
clustered (32 cl.)	5.77	80.56%	80.94%	79.42%	80.13%	82.62%	81.64%	86.32%	82.30%	78.90%	82.60%	69.28%
clustered (64 cl.)	7.03	78.03%	78.00%	77.55%	77.38%	81.48%	79.50%	82.97%	79.90%	76.07%	80.94%	69.66%
clustered (128 cl.)	8.39	81.27%	81.83%	80.19%	80.77%	82.62%	83.09%	86.72%	82.85%	79.82%	83.06%	69.66%
clust. (4 cl.) + quant.	1.39	77.89%	78.76%	76.32%	77.62%	78.49%	78.88%	82.57%	79.56%	76.11%	81.07%	66.66%
clust. (8 cl.) + quant.	2.29	79.97%	80.02%	79.37%	79.46%	82.33%	81.29%	84.80%	81.41%	78.49%	82.33%	68.53%
clust. (16 cl.) + quant.	3.53	76.44%	77.22%	74.99%	75.86%	78.34%	78.32%	80.54%	79.35%	75.23%	77.49%	62.17%
clust. (32 cl.) + quant.	4.71	80.62%	81.13%	79.47%	80.15%	83.04%	81.78%	85.71%	82.57%	79.30%	81.80%	69.28%
clust. (64 cl.) + quant.	5.52	78.00%	78.02%	77.52%	77.37%	80.62%	79.84%	82.57%	79.83%	76.15%	81.07%	68.53%
clust. (128 cl.) + quant.	5.89	81.13%	81.68%	80.09%	80.74%	82.19%	82.53%	86.01%	82.71%	79.76%	83.20%	68.53%