# Efficient Large-Scale Face Clustering Using an Online Mixture of Gaussians

David Montero, Naiara Aginako, Basilio Sierra and Marcos Nieto

*Abstract*— In this work, we address the problem of large-scale online face clustering: given a continuous stream of unknown faces, create a database grouping the incoming faces by their identity. The database must be updated every time a new face arrives. In addition, the solution must be efficient, accurate and scalable. For this purpose, we present an online gaussian mixture-based clustering method (OGMC). The key idea of this method is the proposal that an identity can be represented by more than just one distribution or cluster. Using feature vectors (f-vectors) extracted from the incoming faces, OGMC generates clusters that may be connected to others depending on their proximity and their robustness. Every time a cluster is updated with a new sample, its connections are also updated. With this approach, we reduce the dependency of the clustering process on the order and the size of the incoming data and we are able to deal with complex data distributions. Experimental results show that the proposed approach outperforms state-of-the-art clustering methods on large-scale face clustering benchmarks not only in accuracy, but also in efficiency and scalability.

## I. INTRODUCTION

In recent years, face clustering by identity has been in high demand in a variety of contexts and, in some cases, to support real-time operation of application. For instance, it is required in real-time video-surveillance applications, where there is a need to maintain an updated database of subjects within the monitored area, either to control their locations [1], [2] or to re-identify a subject if necessary [3], [4]. Another usage of face clustering in a real-time application is people flow monitoring in large infrastructures. This type of application aims to generate Performance Indicators, such as "waiting times", "process throughput", "person show-up profile", "queue length overrun" and "area occupancy" [5], [6], which need to be computed for increasingly large number of cameras, and thus demanding a higher level of computation scalability.

An important issue that arises about these types of applications is data privacy [7], [8], [9]. Personal information that could be used to identify the subjects should not be stored (i.e. images of their faces). Recent clustering algorithms rely on Deep Neural Networks (DNN) that can infer feature vectors (f-vectors) from the targeted face images, which correspond to abstract representations of the people appearances used for training. Although in the last years some promising methods for face rendering from f-vectors have been released [10], [11], they require to know the feature

David Montero an Marcos Nieto are with Vicomtech Foundation, Basque Research and Technology Alliance (BRTA), Mikeletegi 57, 20009 Donostia-San Sebastián (Spain). Email: dmontero@vicomtech.org, mnieto@vicomtech.org

Naiara Aginako and Basilio Sierra are with the University of the Basque Country. Email: naiara.aginako@ehu.eus, b.sierra@ehu.eus
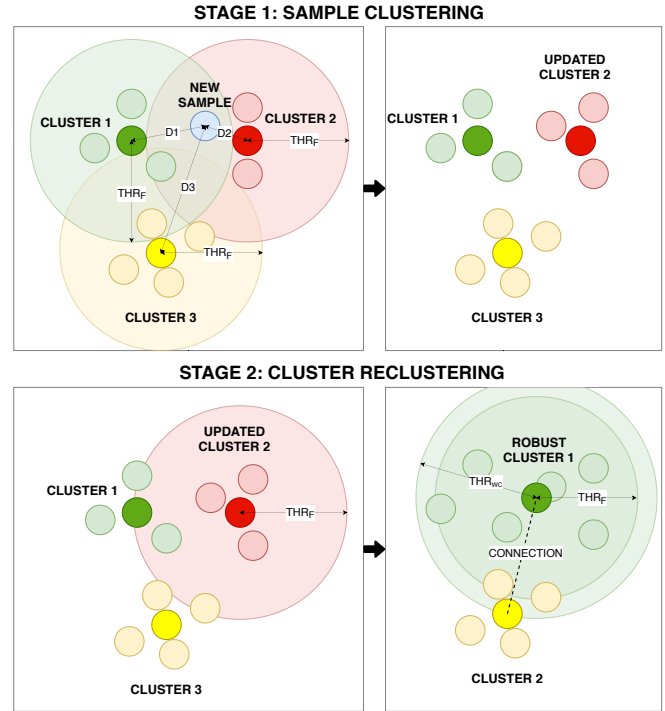
Fig. 1. Example of the operation of the proposed online clustering algorithm when a new sample arrives.

extraction network in order to generate acceptable results [11] or to train the decoding network [10]. Therefore, the individual identities can be protected by hiding and securing the embedding network. Obviously, this approach has limitations, such as the possibility of matching different people that look too similar for the trained model. Therefore, the design and training of the DNN should cover as best as possible the subtle facial appearance dissimilarities. State-of-the-art face recognition models [12], [13] aim at accomplishing that goal.

Furthermore, these real-time applications may work in large-scale unconstrained environments, where no information about the distribution of face representations or the number of identities is available. The faces may come from multiple cameras placed in different locations and positions, so they may have different orientations, lighting conditions, partial occlusions, etc. This will lead to complex data distributions. Most traditional and state-of-the-art clustering methods are offline ([14], [15], [16], [17]). These offline approaches are not suitable for real-time large-scale scenarios, as they need to repeat the whole clustering process

every time a new sample arrives. In addition, most of them have difficulties dealing with complex data distributions.

To overcome these problems, we present an online gaussian mixture-based clustering method (OGMC). Our proposal key idea is that an identity may be represented by several distributions or clusters. Using feature vectors (f-vectors) extracted from the incoming faces, OGMC generates clusters that may be connected to others depending on their proximity and their robustness. Every time a cluster is updated with a new sample, its connections are also updated. With this approach, we reduce the dependency of the clustering process on the order and the size of the incoming data and we are able to deal with complex data distributions. The high-level idea of the method is exposed in figure 1. Experimental results show that our approach outperforms state-of-the-art clustering methods on large-scale face clustering benchmarks not only in accuracy, but also in efficiency and scalability.

The rest of the paper is organized as follows. First, we present a review of the related work in section 2. Section 3 describes the proposed clustering method. In section 4 we provide the experimental results. Finally, the conclusions are given in section 5.

## II. RELATED WORK

### A. Unconstrained Face Clustering

Face clustering in unconstrained environments has become a well-studied topic over the last few years. The addressed scenarios are increasingly complex and encompass a greater number of identities. The huge number of faces and the intra-class appearance changes that might happen due to environmental variations (e.g. pose, illumination, expression, ornaments, occlusions, resolution, image noise) lead to complex distributions of face representations. Traditional clustering algorithms, as K-Means [16] or spectral clustering [18], suffer from this complexity and they are not able to achieve acceptable performance in these scenarios, as they make assumptions on data distribution. For instance, K-Means tends to generate similar-sized clusters.

The new trends combine face recognition models based on DNN to extract facial f-vectors with sophisticated clustering algorithms that can group them in distinguishable identities, despite the intra-class appearance variability. Shi et al. [19] proposed the Conditional Pairwise Clustering (ConPaC) algorithm, which is based on the direct estimation of an adjacency matrix using pairwise similarities between f-vectors. In [14] a linkage-based face clustering algorithm is presented, where a graph convolution network decides which pairs of nodes should be linked. Lin et al. [20] adopted an Agglomerative Hierarchical Clustering approach, considering the distance measure in the embedded space and the dissimilarity between two groups of faces. In [15] an approximate rank-order clustering is presented, which predicts whether a node should be linked to its k Nearest Neighbors (kNN), and transitively merges all linked pairs.

Nevertheless, all these state-of-the-art methodologies employ offline algorithms. They process entirely the gathered data every time a new sample arrives, repeating the clustering process with increasing computational cost, in order to get an optimal result. In addition, most of them suffer from scalability problems in terms of accuracy and processing time. For instance, the complexity of ConPaC can scale up to $O(TN^3)$, where $N$ is the number of f-vectors and $T$ the number of iterations. Wang et al. [14] and Otto et al. [15] reduce the complexity of the proposed algorithms using kNN graphs to reduce the number of comparisons, but the computational cost is still too high to consider them for online applications.

### B. Online Clustering

In recent years, numerous online clustering algorithms have emerged to tackle challenging situations. The main advantage of this type of algorithms is that they are able to process a new sample without repeating the whole clustering process. Therefore, they are the best choice when dealing with large-scale real-time scenarios, but with the added challenge of controlling and defining the learning rate (i.e. how new data updates the learnt models).

These types of algorithms have been applied in a wide variety of contexts. For instance they have gained an increasing importance in text clustering. In [21], the authors proposed an online clustering method for grouping data streams from social networks by their topic, using a similarity measure computed taking into account both the cluster age and the employed terms. Yin and Wang [22] presented an alternative text clustering method, assuming an unknown number of clusters, but below a maximum. Online clustering algorithms have also been employed for unsupervised representation learning [23], where the cluster centroids evolve dynamically, keeping the classifier stably updated. Another example of online clustering application is MalFamAware [24], an algorithm created to group new malwares into families to discern if they are novel or just a variant of a known sample. Even traditional offline methods as K-Means have their own online implementation [25].

Some novel clustering algorithms try to combine both online and offline approaches. Wang and Imura [26] presented a gaussian process-based incremental neural network, where the new samples are treated as nodes, which may be connected to others. When a new sample is clustered, only the connections of its neighbours are re-evaluated. Meanwhile, in [27], samples are grouped in spherical static microclusters, connected together to create dynamically shaped macro-clusters. These approaches aim to achieve the same accuracy as offline methods but with an important decrease in the computation time.

In face clustering there are also some online approaches. In [28], the authors propose an online algorithm for clustering faces in long videos. They process the data sequentially in short segments of variable length and create clusters using face representations and several spatio-temporal constraints. Nevertheless its reliance on these constraints makes their method unsuitable for unconstrained environments and highly susceptible to the order of the incoming data. Tapaswi et al. [29] presented another online face clustering method
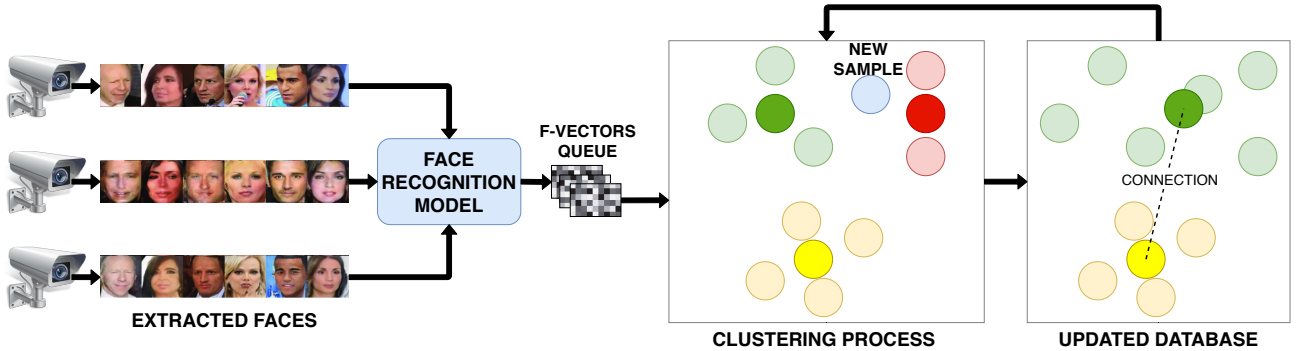
Fig. 2. Example of an online clustering system. A continuous stream of face images, extracted from a set of video-surveillance cameras, are processed by a face recognition model and the extracted f-vectors are enqueued. The online clustering process updates the database with every new sample without repeating the whole process.

for long videos. The algorithm creates spherical clusters with a shared radius which may vary dynamically. They assume that each identity is represented by an identical distribution. Thus, this method is not able to deal with complex data distributions and, therefore, is not a valid solution for unconstrained environments.

Our proposed method aims to cover the need of an online face clustering algorithm capable of working in large-scale unconstrained environments in real time, achieving state-of-the-art results.

## III. PROPOSED METHOD

### A. Problem Definition

We consider the problem of online clustering: given a continuous stream of unknown faces, create a database grouping the incoming faces by their identity. The database must be updated every time a new face arrives, so that information on existing identities is available in real time. An example of an online clustering system is shown in Figure 2. For time and scalability considerations, it is not viable to regenerate the whole database in every iteration, so the algorithm must cluster the new sample using the information from the existing database and update it. Therefore, the problem can be modeled as follows:

$$D_i = F(S_i, D_{i-1})$$
$$D_i = C_i, I_i \tag{1}$$

where $D_i$ is the updated database for the i-th iteration, $F$ is the clustering process repeated in each iteration, $S_i$ is the i-th sample, considering a sample as a normalized N-dimensional f-vector extracted from an incoming face, and $D_{i-1}$ is the resulting database from the previous iteration. The database $D$ is represented by the group of computed clusters $C$ and by the identities $I$ to which they belong. We aim at modeling $F$ maximizing the accuracy and the scalability and minimizing the iteration time.

### B. Expectation-Maximization Approach

The problem modeled by the equation 1 is similar to the one tackled by the Expectation-Maximization (EM) algorithm [30]. The EM algorithm is a well-known iterative

approach to perform maximum likelihood estimation in the presence of latent variables. It has been widely used in clustering applications [31], [32], [33], [34]. The problem formulation is the following: given a statistical model generated by a set of observed data $X$, a set of missing values $Z$, and a vector of unknown parameters $\Theta$, along with a likelihood function:

$$L(\Theta; X, Z) = p(X, Z|\Theta) \tag{2}$$

and the maximum likelihood estimate (MLE) of the unknown parameters is determined by maximizing the marginal likelihood of the observed data:

$$L(\Theta; X) = p(X|\Theta) \tag{3}$$

The EM algorithm aims to find the MLE of the marginal likelihood by iteratively applying two steps:

- Expectation Step: estimates the values of the missing data $Z$, using the observed data $X$ and the current estimation of the parameters $\Theta_i$.

$$Z = F(\Theta_i, X) \tag{4}$$

- Maximization Step: update the parameters $\Theta_{i+1}$ using the observed data $X$ and the new estimated data $Z$.

$$\Theta_{i+1} = F(X, Z) \tag{5}$$

In our context, we apply the EM algorithm assuming that, at each iteration, the observed data is the group of processed samples $S$, the parameters are the features of the computed clusters $C$ and the missing value we want to estimate is the identity of the cluster to which the new sample belongs $I$.

Thus, in the estimation step, using the parameters of the clusters computed with the already processed samples, the algorithm decides whether the new sample should be merged with an existing cluster or create a new one. Then, in the maximization step, the parameters of the involved cluster are updated.

Nevertheless, our method follows a variant of this two-steps approach. Every time a cluster is updated, a new EM algorithm is launched, called cluster reclustering, where we

Fig. 3. Examples of complex identities extracted from the IJB-C dataset. Different variations in lighting conditions, occlusions, perspective and facial attributes such as beard, glasses, hat or hair can be observed. Trying to group all faces in just one cluster may lead to errors due to overly permissive thresholds and poor quality centroids.
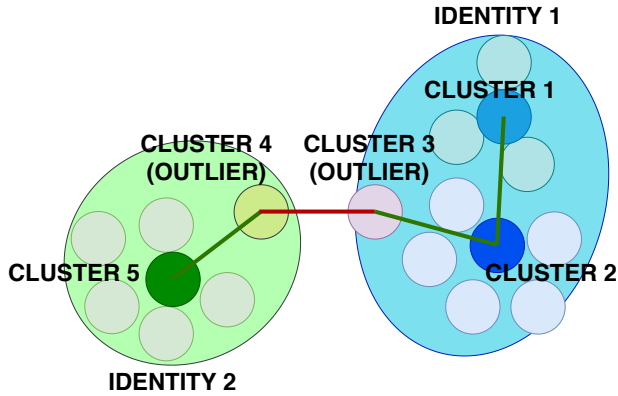


Fig. 4. Example of an erroneous connection caused by outliers. Clusters 1, 2 and 5 are robust while clusters 3 and 4 are non-robust. These problems can be avoided by applying the defined connection rules.

check if the updated cluster can be fused with others. This is an iterative process, so a cluster may be fused multiple times. Thus, our algorithm is divided into two stages:

- Sample clustering stage: the new sample is clustered, updating an existing cluster or creating a new one.
- Cluster reclustering stage: iteratively tries to fuse the updated cluster with the rest of the clusters in the database.

The two-stages process is illustrated in Figure 1, and described in subsequent sections.

### C. Cluster Connections

As we want our algorithm to work in unconstrained environments, it must be able to deal with complex data distributions. Therefore, we must consider the possibility that an identity is represented by more than one cluster. For instance, facial attributes of a person may change for different reasons (i.e. glasses, beard, hair, perspective, lighting conditions, ...), and trying to group all faces in just one cluster may lead to errors due to overly permissive thresholds and poor quality

centroids. Several examples of complex identities extracted from the IJB-C dataset [35] are exposed in Figure 3.

For this reason, we introduce the concept of cluster connection. A connection between two clusters implies that they belong to the same identity, but they represent different distributions in the feature-space. In other words, there is a high similarity between both clusters, but this similarity is not high enough to fuse them (see Figure 1). This way, the algorithm creates trees of connected clusters to deal with complex data distributions.

Furthermore, with these connections we highly reduce the dependency of our algorithm to the order of the incoming samples, as the connections of a cluster are checked and updated every time it is fused with a new sample or with another cluster.

Nevertheless, allowing multiple connections without control may lead to erroneous connections between outliers and clusters belonging to different identities (see Figure 4). To overcome this problem, we create the concept of robust clusters and several rules to control the connections. A cluster is considered robust if it is composed of, at least, $ns_r$ samples, so that we can ensure that it is not a group of outliers. The connection rules are the following:

- A non-robust cluster shall have at most one connection, and it can only be connected to a robust cluster. With this rule, we avoid erroneous connections caused by outliers and redundant connections of non-robust members of an identity, solving problems as the one exposed in Figure 4.
- Two robust clusters can not be fused together. We consider a robust cluster as a valid distribution which represents a subset of an identity samples. Therefore, joining two robust clusters would result in a poorer representation of the identity and a loss of information.
- A robust cluster may have a maximum number of connections $nc_{max}$. This limitation is adopted to reduce the computation time, specially when checking if the connections of a cluster are still valid.

- Every time a cluster is fused or it is connected to new clusters, a process to check connections is triggered. This function checks if the connections of the updated cluster are still valid and that the maximum number of connections is not exceeded. Otherwise, the weakest connections are removed.

### D. Cluster Representation

The next step in the creation of the clustering algorithm is the selection of the group of parameters that represents each cluster. These parameters must contain enough useful information about the cluster they represent in order to obtain accurate results. Furthermore, this information should be brief and concise, as the algorithm must be fast and scalable.

Considering these factors, we decided to model the clusters using multivariate normal distributions, as it only depends on two parameters: the mean vector $\mu$ and the covariance matrix $\Sigma$. This is a design choice also motivated because the EM algorithm works well with these kind of distributions [36], [37], [38]. Therefore, the density function that models the probability that the N-dimensional sample $S_i$ belongs to cluster $j$ is:

$$P(S_i|C_j) \propto exp\left(-\frac{1}{2}\left(S_i - \mu_j{}^T\right)\Sigma_j^{-1}\left(S_i - \mu_j\right)\right) \quad (6)$$

where the Mahalanobis distance [39] between $S_i$ and $\mu_j$ can be directly used to evaluate which is the closer centroid for a certain sample. Indeed, if we assume all dimensions are independent and have the same variance, we can operate on Mahalanobis distances to reduce the computational load of the algorithm, defined as follows:

$$D_M(S_i, \mu_j, \sigma_j) = \frac{1}{\sigma_j} dist(S_i, \mu_j) \quad (7)$$

where $dist(S, \mu) = ||S - \mu||_2$.

The three possible cases when comparing two clusters (fusion, connection or no relation) are then modeled as a mixture of gaussians of these normal distributions, and effectively computed using the Mahalanobis distance. These distributions have all zero mean, which is equal to the minimum possible distance between two centroids. Their deviations $\sigma$ take fixed values depending on whether the compared clusters are robust or not and on the previously defined connection rules (see Figure 5).

Therefore, we define three euclidean distance thresholds to cover all the possible cases generated by the different $\sigma$:

- Fusion threshold $thr_f$: to decide whether two clusters should be fused together.
- Weak connection threshold $thr_{wc}$: to connect a robust cluster with a non-robust one.
- Strong connection threshold $thr_{sc}$: to connect two robust clusters.

The mean of a cluster is represented by its centroid $C^*$, computed by normalizing the sum of the features of all the samples contained in the cluster:
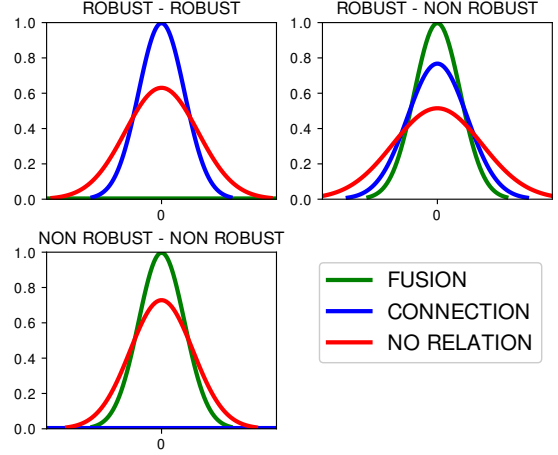


Fig. 5. Normal distributions representing the probability that a cluster can be fused to, connected to or independent of another. The gaussians are centered in 0 (the minimum possible distance between the centroids). The variance depend on whether the compared clusters are robust or not and on the connection rules defined in Section III-C.

$$SC_j = \sum S \in cluster_j$$
$$C_j^* = \frac{SC_j}{||SC_j||_2} \quad (8)$$

Thus, the following information is stored for each cluster:

- N-dimensional centroid ($C^*$)
- Sum of the belonging samples features ($SC$)
- Number of belonging samples ($ns$)
- Index of the belonging samples ($sIdx$)
- Index of the connected clusters ($cIdx$)
- Distance to the connected clusters ($cDist$)

### E. Method Implementation

As described in section III-B, our clustering method is divided into two stages. The first one, the sample clustering stage, is illustrated in Figure 6. In this stage, the new incoming sample is processed. The first step is to compute the distance between the normalized sample vector and the normalized centroids of the existing clusters in the database. Different distances may be considered, but we selected the euclidean distance ($dist$). We made this decision because the employed face recognition model was trained using the cosine similarity [12] and, for normalized vectors, it is inversely proportional to the euclidean distance:

$$dist(V1, V2) = \sqrt{2(1 - cosSim(V1, V2))} \quad (9)$$

The distances are computed in parallel taking advantage of the capabilities of a GPU architecture. This way, we reduce the impact of the clusters number in the processing time. For this reason, the normalized centroids are stored directly in the GPU memory.

Once the distances have been computed, they are copied to the RAM memory and the minimum distance is selected using the CPU. If this distance is less than $thr_f$, the sample is
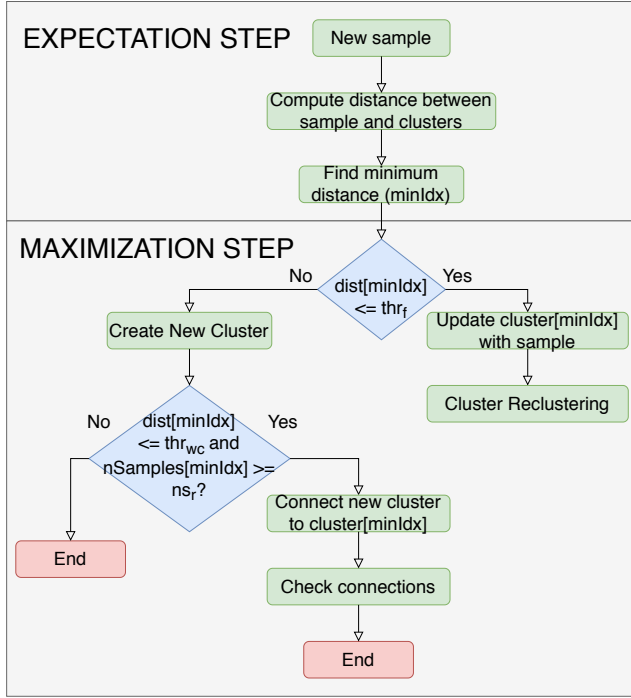
Fig. 6. Diagram describing the sample clustering stage of the proposed clustering algorithm.



Fig. 7. Diagram describing the cluster reclustering stage of the proposed clustering algorithm.

fused with the selected cluster. Otherwise, it is used to create a new cluster. If a new cluster is created, the algorithm checks if it can be connected to the minimum distance cluster. This connection happens if the selected cluster is robust and if the distance is not higher than $thr_{wc}$. If there is a connection, the algorithm checks if the number of connections of the robust cluster has exceeded the maximum permitted ($nc_r$), and, if necessary, erase the weakest connections.

If the new sample $S_i$ is used to update an existing cluster, the cluster reclustering stage is launched (see Figure 7). This stage is composed of an iterative EM algorithm where the updated cluster is tried to be fused or connected with the rest of the clusters in the database, using the parameters $thr_f$, $thr_{wc}$, $thr_{sc}$ and $ns_r$, and the connection rules defined in section III-C.

The selection of the minimum distance is computed in the same way as in the previous stage. If the updated cluster is fused with another one, the distances are computed again for the new updated cluster and the process is repeated. If the updated cluster is connected to another and the updated cluster is robust, the algorithm searches for the next minimum distance and repeats the maximization step.

Before the cluster reclustering stage ends, the algorithm checks if the connections of the updated cluster are still valid, as the centroid of cluster may have changed due to a fusion. Finally, it checks that the number of connection of the updated cluster and of the connected clusters do not exceed the maximum number of connections allowed.

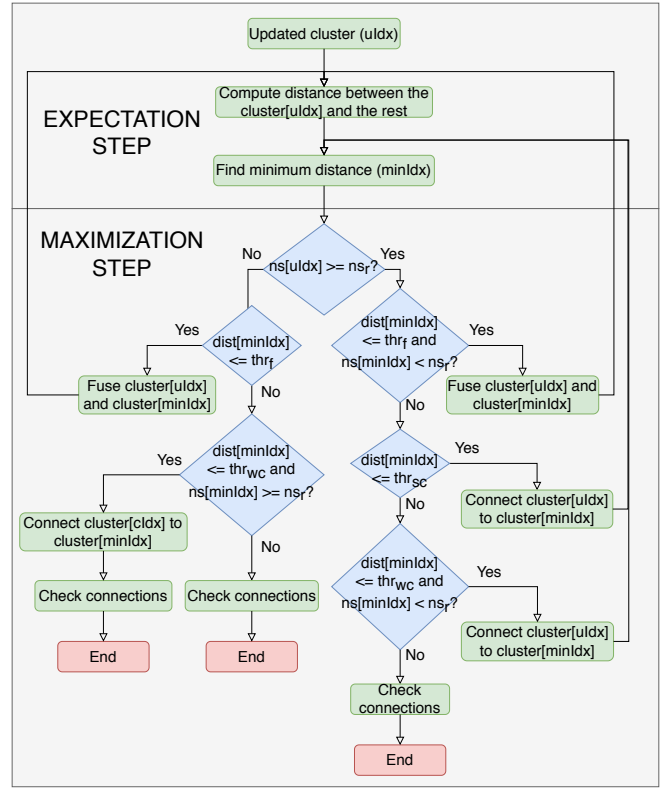The results of the clustering process are extracted after each iteration using a simple recursive function and the vari-ables $sIdx$ and $cIdx$ to merge the samples of all connected clusters.

## IV. EXPERIMENTS

A series of experiments have been conducted to demonstrate the potential of the proposed clustering algorithm OGMC. The experiments are divided into two groups. The first group of aims to measure the performance of OGMC in accuracy and processing time, comparing it with other traditional and state-of-the-art offline clustering methods.

The second group focuses on testing its scalability, measuring the drop in accuracy and the increase in processing time as the number of data samples grows.

Furthermore, an ablation study is presented in order to validate the design decisions, compare the contribution of the different parts of the algorithm and measure the degree of dependency of the model parameters on the face recognition network and the train and test datasets.

Finally, to demonstrate the effectiveness of OGMC beyond face recognition, an additional experiment is conducted with DeepFashion [40], a well-known dataset used for clothes retrieval.

The server used to carry out the experiments was equipped with an NVIDIA Tesla V100 GPU and an Intel Xeon Gold 6230 CPU.

### A. Parameter Tuning

As discussed in section III, the OGMC depends only on 5 parameters: three distance thresholds ($thr_f$, $thr_{wc}$,

$thr_{sc}$), the minimum number of samples to classify a cluster as robust ($ns_r$) and the maximum number of connections allowed for a robust cluster ($nc_r$). This number of parameters is relatively low taking into account that it is an online method which can operate with database regardless of its magnitude, and compared to other state-of-the-art clustering methods as, for example, GCN [14], which requires to train a DNN in addition to three parameters. Furthermore, their values are bounded within certain limits and they are easy to adjust, as it is explained below.

The distance thresholds are floating-point numbers bounded between 0.0 and 2.0, as they represent the euclidean distance between two normalized vectors. These are by far the most sensitive parameters of the algorithm, as small modifications of their values have significant impact on the output.

The minimum number of samples to classify a cluster as robust ($ns_r$) is an integer equal to or greater than 1. On the one hand, the lower its value, the higher number of robust clusters, which may be undesiderable because of errors in connections between outlayers of different identities and an increase in the processing time due to circular connections. On the other hand, if $ns_r$ is too high, the number of robust clusters may not be enough to connect all the non-robust clusters together. We have empirically found that a range between 3 and 6 puts the algorithm into equilibrium, and thus suggest that the optimal value is 4, so the user of the algorithm does not really need to change it.

The maximum number of connections allowed for a robust cluster ($nc_r$) is also an integer equal to or greater than 1. It needs to be adjusted along with $ns_r$, because if $ns_r$ decreases, so does the number of robust clusters and thus increases the number of connections required per robust cluster. If the distance thresholds are set correctly, the accuracy should not be affected by increasing the number of allowed connections. However, this parameter may have a small impact on the processing time by limiting the number of circular connections between clusters of the same identity. We have empirically determined that the appropriate range of values for this parameter is between 5 and 25 and that it can be tuned with a sensitivity of 5.

These parameters, and especially the distance thresholds, depend mainly on the face recognition model, so they only need to be readjusted if the model is replaced. Although they may also have a small dependency on the employed dataset, it has a limited impact on the results, as we report in the subsequent ablation study.

The steps followed to fine tune the parameters are:

- Iterative grid-search for tuning the three distance thresholds. The grid size starts with a size of 0.1 and is halved on each iteration until it reaches a resolution of 0.0125 (4 iterations). The total number of tests is limited according to the restrictions: $thr_{wc}$ must be smaller than $thr_f$ and that $thr_f$ must be smaller than $thr_{sc}$. In this step, $ns_r$ and $nc_r$ take fixed values of 4 and 10 respectively.
- Single grid-search for tuning $ns_r$ and $nc_r$. For $ns_r$ use

a grid size of 1 between 3 and 6 and for $nc_r$ a grid size of 5 between 5 and 25 (20 cells in total). In this step the distance thresholds are fixed and take the values computed in the previous step.

Exploring the parameter space with the suggested approach and testing the response of the algorithm in accuracy and processing time against a training dataset, the user can easily configure the algorithm according to the application requirements. Furthermore, we parallelize every grid search in order to reduce the tuning time.

Finally, the dataset selected for tuning the parameters is reported in the description of each experiment.

### B. Face Clustering Performance

For the first experiment, we use the IJB-B dataset [41], a well-known dataset of unconstrained in-the-wild face images. This dataset includes a clustering protocol consisting of seven subtasks that vary in the number of identities and the number of faces. We select the last subtask, as it is the most challenging one, with the highest number of identities (1,845) and faces (68,195).

For a fair comparison with other methods, and to demonstrate the OGMC algorithm is independent of the recognition model, we use the same vectors as in [14] for the experiment. These vectors have 512 dimensions. For tuning the parameters of OGMC we use the provided features and labels from the CASIA dataset [42]: $thr_f = 1.01$, $thr_{wc} = 1.12$, $thr_{sc} = 0.99$, $ns_r = 5$, $nc_r = 5$.

The performance is measured following the recommendations in [43], selecting the following metrics:

- BCubed F-Measure $F$: represents the clustering system effectiveness, taking the bcubed precision $P$ and recall $R$ into account, which are computed as described in [43].

$$F = 2\frac{P * R}{P + R} \qquad (10)$$

- Normalized Mutual Information (NMI): this measure represents the homogeneity of the clusters. Using the ground truth clusters (G) and the predicted clusters (C) it can be computed with the following equation:

$$NMI(G, C) = \frac{I(G, C)}{\sqrt{H(G)H(C)}} \qquad (11)$$

where $H$ represents the entropy and $I$ is the mutual information.

- Total processing time: since we are comparing our online algorithm with other offline methods, we process all the samples sequentially to simulate an offline behaviour and we compute the time for the whole process.

The results of the experiment are presented in Table I. It can be observed that the proposed method outperforms the others in terms of F-Measure and processing time, while achieving competitive results in the cluster homogeneity measure. Compared to the second best method (GCN-A [14]), OGMC achieves a better F-Measure while reducing the processing time by more than 6 times using the same hardware.

Fig. 8. Example clusters and connections generated by the proposed method in the IJB-C experiment. Each row represents an identity composed of several clusters connected together.

TABLE I

COMPARISON WITH BASELINE METHODS IN TERMS OF BCUBED F-MEASURE, NORMALIZED MUTUAL INFORMATION (NMI) AND PROCESSING TIME USING IJB-B 1845 SUBTASK. SUPERSCRIPT* DENOTES RESULTS REPORTED FROM THE ORIGINAL PAPERS, OTHERWISE ALL METHODS USE THE F-VECTORS FROM [14]. SUPERSCRIPT$^T$ DENOTES TIMES REPORTED FROM [20].

| Method | F-Measure | NMI | Run-Time |
|---|---|---|---|
| K-Means$^T$ [16] | 0.600 | 0.868 | 00:01:00 |
| Spectral$^T$ [18] | 0.516 | 0.785 | - |
| AHC$^T$ [17] | 0.793 | 0.923 | 00:01:32 |
| AP$^T$ [44] | 0.477 | 0.869 | 08:42:50 |
| DBSCAN$^T$ [45] | 0.695 | 0.814 | 00:49:31 |
| ARO$^T$ [15] | 0.755 | 0.913 | 00:01:13 |
| PAHC*$^T$ [20] | 0.610 | 0.890 | 00:03:56 |
| ConPaC*$^T$ [19] | 0.634 | - | 02:53:58 |
| DDC$^T$ [20] | 0.800 | 0.929 | 00:05:32 |
| GCN-A [14] | 0.814 | **0.938** | 00:06:03 |
| **OGMC (ours)** | **0.822** | 0.921 | **00:00:55** |

TABLE II

COMPARISON WITH BASELINE METHODS IN TERMS OF BCUBED F-MEASURE AND PROCESSING TIME USING IJB-C FEATURE VECTORS. ALL METHODS USE THE SAME VECTORS AND HARDWARE.

| Method | F-Measure | Run-Time |
|---|---|---|
| ARO [15] | 0.768 | 00:09:39 |
| GCN [14] | 0.906 | 00:10:32 |
| **OGMC (ours)** | **0.948** | **00:01:32** |

In the second experiment we test the performance of our method using a different face recognition model and a different face dataset. The employed face recognition model is trained using ArcFace loss [12], with ResNet100 [46], [47] as the embedding network and an input resolution of $112 \times 112$. We select MS1MV2 [12] as the training dataset, which is a refinement of MS-Celeb-1M [48].

The dataset selected for this experiment is IJB-C [35], another well-known dataset of unconstrained in-the-wild face images, with a higher number of identities and faces. This dataset also includes a clustering protocol with 8 subtasks. Again, we select the most challenging protocol (IJB-C-3531), with 3,531 identities and 140,623 faces. We use RetinaFace [49] for the face and facial landmarks detection. In order to obtain better quality feature vectors, we filter faces with less than 45 pixels per side and we normalize the face patches applying an affine transformation using reference facial landmarks, as recommendend in [50]. After filtering, 120,661 vectors belonging to 3,529 identities are extracted.

We compare our algorithm with the two best state-of-the-art methods: GCN [14], which achieved the highest accuracy in the first experiment (without considering ours), and ARO [15], which achieved the best trade-off between accuracy and speed (without considering ours). We adjust the parameters of both methods to achieve the best performance. Furthermore, for GCN, we retrain the network using a subset of VGG2 dataset [51], containing over 300k images and 8500 identities, during 4 epochs (following the recommendations in [14]). Finally, we readjust the parameters of our algorithm, as we are using a different face recognition model, using the same subset of VGG2 employed for retraining GCN: $thr_f$ = 1.07, $thr_{wc}$ = 1.15, $thr_{sc}$ = 1.05, $ns_r$ = 5, $nc_r$ = 10.

The results of this experiment, presented in Table II, show that our method outperform the others in terms of F-Measure and processing time. OGMC runs 7 times faster than GCN and more than 6 times faster than ARO. These ratios of processing time show that our method is also more suitable to scale-up compared to the others.

We also present several example clusters generated by the proposed method during the experiment. In Figure 8, each row represents an identity composed of different clusters connected to each other. These identities contain variations in lighting conditions, partial occlusions, perspective and facial attributes, so they are represented by complex data

distributions. With our proposed method, we are able to accurately approximate these distributions using a variable number of connected clusters. Each robust cluster generates a centroid (representative f-vector) which represents the identity under certain condition range. For instance, in the second row of Figure 8, the second cluster centroid represents the identity under dim lighting conditions. If we tried to group all faces of this identity in just one cluster, it could lead to errors due to overly permissive thresholds and worse quality centroids. Another benefit of OGMC that can be observed in this figure is that the identity outliers are contained in non-robust clusters connected to the robust ones. This way, the quality of the robust cluster centroids is not compromised by these outliers and the number of matching errors are reduced. A clear example of this behaviour is exposed in the last row of Figure 8, where the second and the third clusters are outliers generated by a combination of unsuitable conditions (lighting, blurring, occlusion...).

In the last experiment of this section we test the accuracy of OGMC under extreme conditions: using a face recognition model which extracts vectors with less features and using a database with a much higher number of identities. Again, as we want a fair comparison with other methods we decide to replicate the experiment proposed in [52]. In this experiment, they selected a subset of the database MS-Celeb-1M [48] that contains 5.8M images from 86K identities and randomly split it into 10 parts with an almost equal number of identities. Then, they randomly selected 1 part as labeled data for training and the other 9 parts as unlabeled data. With the unlabeled data, they created 5 tests with an increasing number of vectors and identities. The last test has 5.2M vectors and 77K identities. Furthermore the provided vectors have only 256 features, compared to the 512 previously used.

Thus, as they did with the rest of the methods, we tune the parameters of OGMC using the provided training data: $thr_f = 0.85$, $thr_{wc} = 1.02$, $thr_{sc} = 0.72$, $ns_r = 4$, $nc_r = 25$. Then, we evaluate the algorithm using the 5 test subsets with an increasing scale of vectors and identities. The results of the experiment are presented in Table III. It can be observed that OGMC outperforms the rest of the methods consistently in every test subset.

With this experiment, we also prove that the parameters of OGMC do not need to be readjusted if the scale of the dataset increases, so the user just need to tune the parameters if the face recognition model is changed, as discussed above. Of course, as for any other method, we must ensure the training data is large and varied enough to extract the necessary information from the model to correctly tune the parameters.

### C. Clustering Scalability

We also want to prove that our method is scalable and that it is suitable for large-scale real-time applications. For this reason, we conduct the following two experiments. In the first one, we test the scalability of our method adding 1, 2 and 3 millions of distractors to the IJB-C vectors used in the previous experiment. Thus, we can observe how the accuracy and the processing time evolve with the size of

TABLE III

COMPARISON WITH BASELINE METHODS IN TERMS OF BCUBED F-MEASURE USING SUBSETS OF DIFFERENT SIZES FROM MS-CELEB-1M DATASET. ALL METHODS USE THE SAME 256-DIMENSIONAL VECTORS PROVIDED BY [52].

| Method | Test | | | | |
| | Number of samples | | | | |
| | 584K | 1.74M | 2.89M | 4.05M | 5.21M |
| | Number of identities | | | | |
| | 8.5K | 25.7K | 42.8K | 60.0K | 77.1K |
| | BCubed F-Measure | | | | |
|---|---|---|---|---|---|
| K-means [16], [53] | 0.812 | 0.752 | 0.723 | 0.706 | 0.694 |
| HAC [54] | 0.705 | 0.695 | 0.686 | 0.677 | 0.670 |
| DBSCAN [45] | 0.672 | 0.665 | 0.663 | 0.449 | 0.447 |
| ARO [15] | 0.170 | 0.124 | 0.110 | 0.105 | 0.100 |
| CDP [55] | 0.787 | 0.758 | 0.746 | 0.736 | 0.729 |
| GCN [14] | 0.844 | 0.816 | 0.801 | 0.793 | 0.786 |
| LTC [56] | 0.855 | 0.830 | 0.811 | 0.798 | 0.789 |
| GCN-V [52] | 0.858 | 0.826 | 0.811 | 0.799 | 0.791 |
| GCN-(V+E) [52] | 0.861 | 0.828 | 0.812 | 0.801 | 0.793 |
| **OGMC (ours)** | **0.906** | **0.881** | **0.864** | **0.851** | **0.839** |

TABLE IV

RESULTS OF THE PROPOSED METHOD ON IJB-C EXPERIMENT ADDING DISTRACTORS FROM VGG2 DATASET.

| Samples Number | F-Measure | Run-Time |
|---|---|---|
| IJB-C | 0.952 | 00:01:32 |
| IJB-C + 1 million | 0.951 | 00:24:55 |
| IJB-C + 2 millions | 0.948 | 01:00:40 |
| IJB-C + 3 millions | 0.947 | 02:13:02 |

the database. We generate the distractors from the VGG2 face dataset [51], which contains 3.3 millions faces belonging to more than 9000 identities, with large variations in pose, age, illumination, ethnicity and profession. The same face recognition model of the previous IJB-C experiment is used, so we use the same parameter values: $thr_f = 1.04$, $thr_{wc} = 1.13$, $thr_{sc} = 1.03$, $ns_r = 5$, $nc_r = 5$.

Tests are executed 10 times, shuffling all the vectors a random number of times and computing the average F-Measure and processing time for all tests. For the computation of the F-Measure the distractors are ignored. The results shown in Table IV, demonstrate that OGMC has an extremelly high scalability, with a drop in accuracy of 0.1% when adding 1 million of distractors and 0.5% when adding 3 millions. Furthermore, OGMC able to process more than 1.1 million faces in less than 25 minutes and more than 3.1 million faces in 2 hours and 13 minutes. Some examples of the clusters and connections generated in this experiment are shown in Figure 11.

The last experiment aims to demonstrate that OGMC is suitable for real-time applications. It consist of repeating the previous test using IJB-C vectors with 3 millions of distractors and measuring how the processing time per sample evolves with the number of processed samples and with the number of clusters in the database (note that the number of clusters is not equal to the number of identities, as we are
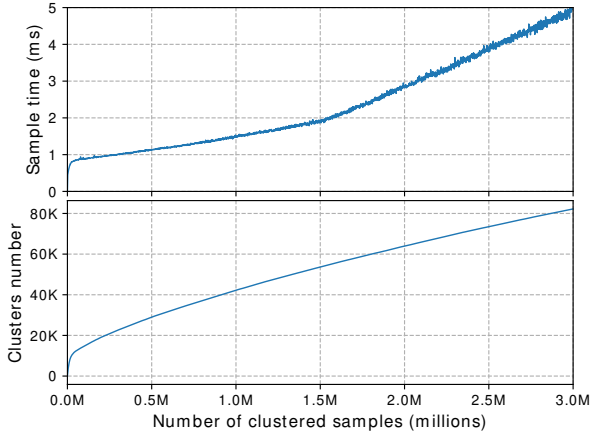
Fig. 9. Evolution of the clustering time per sample with the number of samples clusterized compared to the evolution of the number of clusters in the database.

not taking connections into account).

As in the previous experiment, the test is repeated 10 times, shuffling all the vectors a random number of times and computing the results as the average value of all tests. Figure 9 shows that, after processing 3 millions samples and with a database of 80 thousand clusters, the processing time per sample is still 5 milliseconds, which still allows OGMC to process 200 samples per second in real-time. This results demonstrate that the proposed method is suitable for real-time applications, even when dealing with extremely large amounts of data.

*D. Ablation Study*

We present an ablation study to discuss about the impact of design choices on the performance of OGMC. First, we demonstrate the contribution of the reclustering stage of the algorithm. For this purpose, we run the IJB-C with VGG2 distractors experiment applying only the first stage of OGMC and compare the results with the ones obtained applying the full method. Table V shows the results of the experiment. Comparing the processing times, we can see that removing the reclustering stage significantly reduces the complexity of the algorithm. However, it also leads to an important drop in accuracy ($\approx$4%). Furthermore, removing the second stage also significantly increases the dependency of accuracy on the order of the input data. This is because the connections are not re-evaluated when the clusters are updated.

We also want to analyze the speed up added by the GPU parallelization. We reimplement the distance computing module with CPU and repeat the IJB-C experiment with both versions of OGMC to measure how the processing time evolves with the number of clusterized samples in each case. The results of the comparison in Figure 10 show that using the GPU for the distance computing parallelization produces a huge decrease on the computation time per sample and a huge increase on the scalability. However, this simple task only consumes a small percentage of the GPU utilization

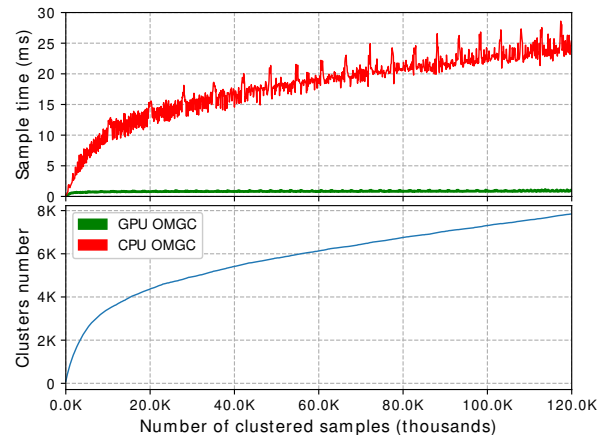| Dataset | Full method | | First Stage | |
|---|---|---|---|---|
| | F-Meas | Run-Time | F-Meas | Run-Time |
| IJB-C | 0.952 | 00:01:32 | 0.914 | 00:00:49 |
| IJB-C + 1 million | 0.951 | 00:24:55 | 0.912 | 00:10:12 |
| IJB-C + 2 millions | 0.948 | 01:00:40 | 0.907 | 00:23:11 |
| IJB-C + 3 millions | 0.947 | 02:13:02 | 0.901 | 00:59:18 |



Fig. 10. Comparison of the evolution of the clustering time per sample with the number of samples clusterized for GPU and CPU versions of OGMC.

compared to other methods based on neural networks as GCN [14] or GCN-V [52]. Furthermore the GPU memory needed for the cluster centroids is very limited (2.4GB for 1 million 512-dimensional vectors) and can be allocated dynamically as the database growths, which makes OGMC suitable to run together with other GPU-based algorithms (for example a DNN-based face recognition model).

Finally, we want to demonstrate that the parameters of OGMC do not have a significant dependency on the training dataset, as long as it is large and varied enough to extract the necessary information from the face recognition model. In order to prove this, we repeat the IJB-C experiment with the same face recognition model but, instead of using the VGG2 subset, we tune the parameters using the IJB-C vectors, so that we obtain the best possible results in this test. Thus, we obtain the following values for the parameters: $thr_f$ = 1.04, $thr_{wc}$ = 1.13, $thr_{sc}$ = 1.03, $ns_r$ = 5, $nc_r$ = 5. With these values OGMC achieves a BCubed F-Measure of 0.952, only a 0.4% of improvement over the original test. Thus, this experiment suggests our hypothesis is correct, the OGMC parameters depend almost entirely on the face recognition model and they only need to be readjusted if such model is replaced.

*E. Beyond Face Recogniton*

To conclude the experimental section, we evaluate the effectiveness of OGMC for tasks beyond face recognition.

TABLE VI

COMPARISON WITH BASELINE METHODS IN TERMS OF BCUBED
F-MEASURE WITH DEEPFASHION DATASET. ALL METHODS USE THE
SAME VECTORS FROM [52].

| Method | F-Measure |
|---|---|
| K-Means [16], [53] | 0.538 |
| HAC [54] | 0.488 |
| DBSCAN [45] | 0.532 |
| MeanShift [57], [58] | 0.567 |
| Spectral [59], [60] | 0.464 |
| ARO [15] | 0.530 |
| CDP [55] | 0.578 |
| GCN [14] | 0.589 |
| LTC [56] | 0.591 |
| GCN-V [52] | 0.573 |
| GCN-(V+E) [52] | 0.601 |
| **OGMC (ours)** | **0.620** |

For this purpose, we reproduce the experiment presented in [52] with DeepFashion dataset [40], a well-known dataset for clothes retrieval. In this experiment they mix the training and testing features in the original split, and randomly sample 25,752 images from 3,997 categories for training and the other 26,960 images with 3,984 categories for testing. For a fair comparison, we use the same 256-dimensional vectors provided in their official repository. Thus, we tune the parameters of OGMC using the training set: $thr_f = 0.51$, $thr_{wc} = 0.58$, $thr_{sc} = 0.49$, $ns_r = 4$, $nc_r = 10$. The results of the experiment are presented in Table VI. Among all the tested methods, OGMC achieves the best F-Measure, demonstrating its suitability for tasks beyond face recognition.

## V. CONCLUSIONS

In this work, we address the problem of large-scale online face clustering. We propose a method based on EM algorithm and Mixture of Gaussians. We introduce the concept of cluster connection, where an identity can be represented by multiple clusters. With this approach, we reduce the dependency of the clustering process on the order and the size of the incoming data and we are able to deal with complex data distributions. The conducted experiments and their derived results show that our method outperforms the state-of-the-art clustering methods, regardless of whether they are online or offline, in terms of accuracy, processing time and scalability.

Future work will focus on further reducing the processing time of the proposed method by also parallelizing the minimum distance search using the computational capabilities of the GPU. In addition, given the current situation with Covid-19, we will study the impact of the use of surgical masks on the clustering process. We believe that for an identity, the proposed method would group faces with and without masks in different clusters, but connected to each other, since we have proved that the algorithm works well with partial occlusions. Finally, as we believe in the open source community, we will soon release the full code of OGMC and all the supplementary material used for the experiments, so anyone can replicate them and contribute to improve the method.

## REFERENCES

[1] Z. Lei, C. Wang, Q. Wang, and Y. Huang, "Real-time face detection and recognition for video surveillance applications," in *2009 WRI World Congress on Computer Science and Information Engineering*, vol. 5, 2009, pp. 168–172.

[2] F. Mahdi, M. Habib, A. Moslehuddin, P. Vasant, S. Mckeever, and M. A. R. Ahad, "Face recognition-based real-time system for surveillance," *Intelligent Decision Technologies*, vol. 11, pp. 1–14, 09 2016.

[3] P. Apoorva., H. C. Impana., S. L. Siri., M. R. Varshitha., and B. Ramesh., "Automated criminal identification by face recognition using open computer vision classifiers," in *2019 3rd International Conference on Computing Methodologies and Communication (ICCMC)*, 2019, pp. 775–778.

[4] W. Yimyam, T. Pinthong, N. Chumuang, and M. Ketcham, "Face detection criminals through cctv cameras," in *2018 14th International Conference on Signal-Image Technology Internet-Based Systems (SITIS)*, 2018, pp. 351–357.

[5] C. A. Mayer, R. Felkel, and K. Peterson, "Best practice on automated passenger flow measurement solutions," in *Journal of Airport Management*, vol. 9, 2015, pp. 144–153.

[6] K.-M. Chien, T.-C. Wu, and T. Luor, "Face recognition and smart people-counting system: Cases of asian trade shows," *Journal of Internet Technology*, vol. 20, no. 2, pp. 435–446, 2019. [Online]. Available: https://jit.ndhu.edu.tw/article/view/2017

[7] E. M. Newton, L. Sweeney, and B. Malin, "Preserving privacy by de-identifying face images," *IEEE Transactions on Knowledge and Data Engineering*, vol. 17, no. 2, pp. 232–243, 2005.

[8] K. W. Bowyer, "Face recognition technology: security versus privacy," *IEEE Technology and Society Magazine*, vol. 23, no. 1, pp. 9–19, 2004.

[9] Z. Erkin, M. Franz, J. Guajardo, S. Katzenbeisser, I. Lagendijk, and T. Toft, "Privacy-preserving face recognition," in *Privacy Enhancing Technologies*, I. Goldberg and M. J. Atallah, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2009, pp. 235–253.

[10] S. Lombardi, J. M. Saragih, T. Simon, and Y. Sheikh, "Deep appearance models for face rendering," *ACM Transactions on Graphics (TOG)*, vol. 37, pp. 1 – 13, 2018.

[11] C. N. Duong, T.-D. Truong, K. G. Quach, H. Bui, K. Roy, and K. Luu, "Vec2face: Unveil human faces from their blackbox features in face recognition," *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6131–6140, 2020.

[12] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "Arcface: Additive angular margin loss for deep face recognition," in *2019 IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, June 2019, pp. 4685–4694.

[13] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 815–823.

[14] Z. Wang, L. Zheng, Y. Li, and S. Wang, "Linkage based face clustering via graph convolution network," *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1117–1125, 2019.

[15] C. Otto, D. Wang, and A. K. Jain, "Clustering millions of faces by identity," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 2, pp. 289–303, Feb 2018.

[16] S. P. Lloyd, "Least squares quantization in pcm," *IEEE Trans. Inf. Theory*, vol. 28, pp. 129–136, 1982.

[17] A. K. Jain and R. C. Dubes, *Algorithms for Clustering Data*. USA: Prentice-Hall, Inc., 1988.

[18] J. Shi and J. Malik, "Normalized cuts and image segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 22, pp. 888–905, 2000.

[19] Y. Shi, C. Otto, and A. K. Jain, "Face clustering: Representation and pairwise constraints," *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 7, pp. 1626–1640, July 2018.

[20] W. Lin, J. Chen, C. D. Castillo, and R. Chellappa, "Deep density clustering of unconstrained faces," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8128–8137.

[21] C. Comito, C. Pizzuti, and N. Procopio, "Online clustering for topic detection in social data streams," in *2016 IEEE 28th International Conference on Tools with Artificial Intelligence (ICTAI)*, 2016, pp. 362–369.

[22] J. Yin and J. Wang, "A text clustering algorithm using an online clustering scheme for initialization," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ser. KDD '16.  New York, NY, USA: Association for Computing Machinery, 2016, p. 1995–2004. [Online]. Available: https://doi.org/10.1145/2939672.2939841

[23] X. Zhan, J. Xie, Z. Liu, Y.-S. Ong, and C. C. Loy, "Online deep clustering for unsupervised representation learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[24] G. Pitolli, G. Laurenza, L. Aniello, L. Querzoni, and R. Baldoni, "Malfamaware: automatic family identification and malware classification through online clustering," *International Journal of Information Security*, 06 2020.

[25] E. Liberty, R. Sriharsha, and M. Sviridenko, "An algorithm for online k-means clustering," in *ALENEX*, 2016.

[26] X. Wang and J. Imura, "A gaussian process-based incremental neural network for online clustering," in *2019 IEEE International Conference on Smart Cloud (SmartCloud)*, 2019, pp. 143–148.

[27] R. Hyde, P. Angelov, and A. Mackenzie, "Fully online clustering of evolving data streams into arbitrarily shaped clusters," *Information Sciences*, vol. 382, pp. 1–41, 12 2016.

[28] P. Kulshreshtha and T. Guha, "An online algorithm for constrained face clustering in videos," in *2018 25th IEEE International Conference on Image Processing (ICIP)*, 2018, pp. 2670–2674.

[29] M. Tapaswi, M. T. Law, and S. Fidler, "Video face clustering with unknown number of clusters," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.

[30] T. K. Moon, "The expectation-maximization algorithm," *IEEE Signal Processing Magazine*, vol. 13, no. 6, pp. 47–60, 1996.

[31] S. Hwang, J. Oh, J. Cox, S. J. T. M.D., and H. F. T. M.D., "Blood detection in wireless capsule endoscopy using expectation maximization clustering," in *Medical Imaging 2006: Image Processing*, J. M. Reinhardt and J. P. W. Pluim, Eds., vol. 6144, International Society for Optics and Photonics.  SPIE, 2006, pp. 577 – 587. [Online]. Available: https://doi.org/10.1117/12.654109

[32] Y. G. Jung, M. S. Kang, and J. Heo, "Clustering performance comparison using k-means and expectation maximization algorithms," *Biotechnology & Biotechnological Equipment*, vol. 28, no. sup1, pp. S44–S48, 2014, pMID: 26019610. [Online]. Available: https://doi.org/10.1080/13102818.2014.949045

[33] S. Nasser, R. Alkhaldi, and G. Vert, "A modified fuzzy k-means clustering using expectation maximization," in *2006 IEEE International Conference on Fuzzy Systems*, 2006, pp. 231–235.

[34] J. Garriga, J. Palmer, A. Oltra, and F. Bartumeus, "Expectation-maximization binary clustering for behavioural annotation," *Mov Ecol*, vol. 11, 03 2015.

[35] B. Maze, J. Adams, J. A. Duncan, N. Kalka, T. Miller, C. Otto, A. K. Jain, W. T. Niggel, J. Anderson, J. Cheney, and P. Grother, "Iarpa janus benchmark - c: Face dataset and protocol," in *2018 International Conference on Biometrics (ICB)*, Feb 2018, pp. 158–165.

[36] J. P. Vila and P. Schniter, "Expectation-maximization gaussian-mixture approximate message passing," *IEEE Transactions on Signal Processing*, vol. 61, no. 19, pp. 4658–4672, 2013.

[37] G. Tian, Y. Xia, Y. Zhang, and D. Feng, "Hybrid genetic and variational expectation-maximization algorithm for gaussian-mixture-model-based brain mr image segmentation," *IEEE Transactions on Information Technology in Biomedicine*, vol. 15, no. 3, pp. 373–380, 2011.

[38] C. Guo, H. Fu, and W. Luk, "A fully-pipelined expectation-maximization engine for gaussian mixture models," in *2012 International Conference on Field-Programmable Technology*, 2012, pp. 182–189.

[39] R. De Maesschalck, D. Jouan-Rimbaud, and D. Massart, "The mahalanobis distance," *Chemometrics and Intelligent Laboratory Systems*, vol. 50, pp. 1–18, 01 2000.

[40] Z. Liu, P. Luo, S. Qiu, X. Wang, and X. Tang, "Deepfashion: Powering robust clothes recognition and retrieval with rich annotations," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 1096–1104.

[41] C. Whitelam, E. Taborsky, A. Blanton, B. Maze, J. Adams, T. Miller, N. Kalka, A. K. Jain, J. A. Duncan, K. Allen, J. Cheney, and P. Grother, "Iarpa janus benchmark-b face dataset," in *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, July 2017, pp. 592–600.

[42] D. Yi, Z. Lei, S. Liao, and S. Z. Li, "Learning face representation from scratch," 2014.

[43] E. Amigó, J. Gonzalo, J. Artiles, and M. Verdejo, "Amigó e, gonzalo j, artiles j et ala comparison of extrinsic clustering evaluation metrics based on formal constraints. inform retriev 12:461-486," *Information Retrieval*, vol. 12, pp. 461–486, 10 2009.

[44] B. Frey and D. Dueck, "Clustering by passing messages between data points," *Science (New York, N.Y.)*, vol. 315, pp. 972–6, 03 2007.

[45] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu, "A density-based algorithm for discovering clusters in large spatial databases with noise," in *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, ser. KDD'96.  AAAI Press, 1996, p. 226–231.

[46] D. Han, J. Kim, and J. Kim, "Deep pyramidal residual networks," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 6307–6315.

[47] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016, pp. 770–778.

[48] Y. Guo, L. Zhang, Y. Hu, X. He, and J. Gao, "Ms-celeb-1m: A dataset and benchmark for large-scale face recognition," in *ECCV*, vol. 9907, 10 2016, pp. 87–102.

[49] J. Deng, J. Guo, E. Ververas, I. Kotsia, and S. Zafeiriou, "Retinaface: Single-shot multi-level face localisation in the wild," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[50] H. Wang, Y. Wang, Z. Zhou, X. Ji, D. Gong, J. Zhou, Z. Li, and W. Liu, "Cosface: Large margin cosine loss for deep face recognition," in *2018 IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, June 2018, pp. 5265–5274.

[51] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman, "Vggface2: A dataset for recognising faces across pose and age," in *2018 13th IEEE International Conference on Automatic Face Gesture Recognition (FG 2018)*, 2018, pp. 67–74.

[52] L. Yang, D. Chen, X. Zhan, R. Zhao, C. C. Loy, and D. Lin, "Learning to cluster faces via confidence and connectivity estimation," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 13 366–13 375.

[53] D. Sculley, "Web-scale k-means clustering," in *Proceedings of the 19th International Conference on World Wide Web*, ser. WWW '10. New York, NY, USA: Association for Computing Machinery, 2010, p. 1177–1178. [Online]. Available: https://doi.org/10.1145/1772690.1772862

[54] R. Sibson, "SLINK: An optimally efficient algorithm for the single-link cluster method," *The Computer Journal*, vol. 16, no. 1, pp. 30–34, 01 1973. [Online]. Available: https://doi.org/10.1093/comjnl/16.1.30

[55] X. Zhan, Z. Liu, J. Yan, D. Lin, and C. C. Loy, "Consensus-driven propagation in massive unlabeled data for face recognition," in *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018.

[56] L. Yang, X. Zhan, D. Chen, J. Yan, C. C. Loy, and D. Lin, "Learning to cluster faces on an affinity graph," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 2293–2301.

[57] Yizong Cheng, "Mean shift, mode seeking, and clustering," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 17, no. 8, pp. 790–799, 1995.

[58] D. Comaniciu and P. Meer, "Mean shift analysis and applications," in *Proceedings of the Seventh IEEE International Conference on Computer Vision*, vol. 2, 1999, pp. 1197–1203 vol.2.

[59] J. Ho, Ming-Husang Yang, Jongwoo Lim, Kuang-Chih Lee, and D. Kriegman, "Clustering appearances of objects under varying illumination conditions," in *2003 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2003. Proceedings.*, vol. 1, 2003, pp. I–I.

[60] A. Y. Ng, M. I. Jordan, and Y. Weiss, "On spectral clustering: Analysis and an algorithm," in *Proceedings of the 14th International Conference on Neural Information Processing Systems: Natural and Synthetic*, ser. NIPS'01.  Cambridge, MA, USA: MIT Press, 2001, p. 849–856.

**David Montero** received his B.Sc. degree in Industrial Engineering and his M.S. degree in Industrial Engineering from ETSI of the University of Seville, Spain, in 2015 and 2017 respectively. Since 2017 he has worked as a Researcher on computer vision and machine learning techniques of the Intelligent Transport Systems and Engineering Area of Vicomtech. Since 2019 he is also working towards the Ph.D. degree with the Department of Computer Science and Artificial Intelligence of the University of the Basque Country. His research interests include pattern recognition and computer vision.

**Naiara Aginako** received her degree in Telecommunications Engineering from the University of the Basque Country (UPV/EHU) in 2005, and her Ph.D. degree in computer sciences from the University of the Basque Country in 2017. From 2003 to 2005 she collaborated with the Signal Processing and Communication Group of the Department of Electronics and Telecommunications. She worked as a senior researcher in Vicomtech from 2005 until 2015, in the Digital Media department. From 2015 until 2018 she taught in the Engineering School of Gipuzkoa in the Applied Mathematics department. Since 2018 she has been working in the Faculty of Informatics in the Computer Science and Artificial Intelligence department (UPV/EHU). Her research is mainly based on Computer Vision more specifically on Image and Video Analysis. Her work as a researcher includes a series of publications and two patents.

**Basilio Sierra** received the B.Sc. degree in computer sciences, the M.Sc. degree in computer science and architecture, and the Ph.D. degree in computer sciences from the University of the Basque Country, in 1990, 1992, and 2000, respectively. He is currently a Full Professor with the Computer Sciences and Artificial Intelligence Department, University of the Basque Country. He is also the Co-Director of the Robotics and Autonomous Systems Group, Donostia-San Sebastian. He is also a Researcher in the fields of robotics and machine learning, where he is working on the use of different paradigms to improve behaviours. He has published more than 50 journal articles, and several book chapters and conference papers.

**Marcos Nieto** received his M.S. and Ph.D. degree in Electrical Engineering from ETSIT of the Universidad Politécnica de Madrid (UPM), Spain, in 2005 and 2010, respectively. From 2005 to 2010 he worked as Researcher within the Image Processing Group at the UPM. Since 2010 he has worked as a Researcher on computer vision and machine learning techniques of the Intelligent Transport Systems and Engineering Area of Vicomtech. He has been the technical and scientific coordinator of FP7 and H2020 projects and has experience in transferring technology to industry to support innovation. He is the author of more than 60 peer reviewed international publications in relevant conferences (+40) and journals (+15). He is also an active reviewer of prestigious journals for IEEE, Elsevier and Springer.

Fig. 11. Example clusters and connections generated by the proposed method in the IJB-C + 3 millions of distractors from VGG2 experiment. Each group of contiguous faces represents a cluster and each line represents a connection.