

Inferring Dynamic Representations of Facial Actions from a Still Image

Siyang Song¹, Enrique Sánchez-Lozano¹, Linlin Shen², Alan Johnston³ and Michel Valstar¹

¹ School of Computer Science, University of Nottingham, UK

² College of Computer Science and Software Engineering, Shenzhen University, China

³ School of Psychology, University of Nottingham, UK

Abstract

Facial actions are spatio-temporal signals by nature, and therefore their modeling is crucially dependent on the availability of temporal information. In this paper, we focus on inferring such temporal dynamics of facial actions when no explicit temporal information is available, i.e. from still images. We present a novel approach to capture multiple scales of such temporal dynamics, with an application to facial Action Unit (AU) intensity estimation and dimensional affect estimation. In particular, 1) we propose a framework that infers a dynamic representation (DR) from a still image, which captures the bi-directional flow of time within a short time-window centered at the input image; 2) we show that we can train our method without the need of explicitly generating target representations, allowing the network to represent dynamics more broadly; and 3) we propose to apply a multiple temporal scale approach that infers DRs for different window lengths (MDR) from a still image. We empirically validate the value of our approach on the task of frame ranking, and show how our proposed MDR attains state of the art results on BP4D for AU intensity estimation and on SEMAINE for dimensional affect estimation, using only still images at test time.

1. Introduction

Temporal dynamics are an important source of information for video-based face analysis. In recent years, many methods have been proposed to exploit that for tasks where the temporal information correlates over time with the target signals [6, 16, 51, 48]. The temporal modeling can be accomplished either by generating a single set of features from multiple consecutive frames at a time (early modeling of temporal dynamics [51]) or by using memory-based models, such as Recurrent Neural Networks (RNNs) or Markov Models (late modeling of temporal dynamics [6]), or by a combination of both [16]. Recently, there is an increasing interest in the early modeling of temporal dynamics, as these can be straightforwardly used in simple CNN-

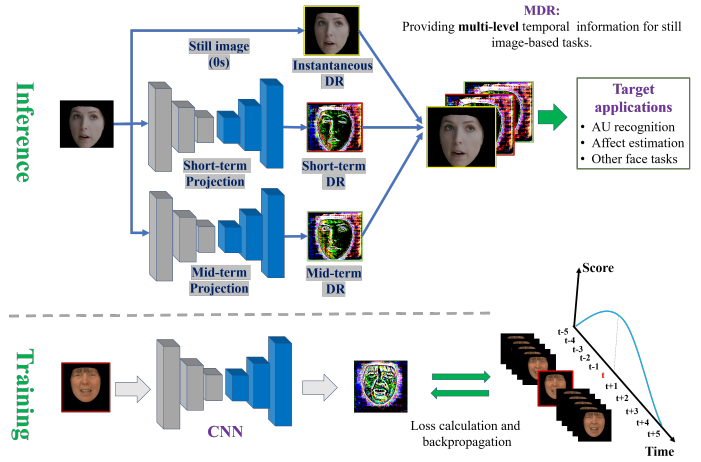


Figure 1. We propose a novel approach to the modeling of temporal face dynamics from still images (top). Our approach can be used to infer several time-length dynamics, and further combined to enhance face-related tasks such as AU intensity estimation and dimensional affect estimation. During training (bottom), a set of videos is used to learn the DRs, without explicitly generating a fixed set of target representations.

based architectures. Among these methods, the dynamic image [1] has achieved promising results at the task of summarizing short-term motion in a fixed-size representation, for the task of action recognition.

While the use of temporal dynamics is desirable for facial behaviour-related tasks, it is often not available. For example, in many applications across domains ranging from medicine to marketing, the analysis of facial expressions or emotions is required from still images. Without temporal information, the performance of state of the art methods for facial expression recognition or affect estimation degrade substantially [21, 16]. While some works have been proposed to anticipate or predict the next frame from a still image or from a video [30, 45], none of existing works have attempted to model the dynamics of facial expressions from a still image.

In this paper, we want to generate a dynamic representa-

tion that summarizes motion, but that can be inferred from still images. To the best of our knowledge, this is the first work that brings the advantages of DRs at summarizing sequences to scenarios where only still images are given. Even though some early works exist that summarize temporal dynamics from image sequences [1], or attempt to predict or anticipate motion from still images [30, 45], to the best of our knowledge, there are no works trying to infer a DR from a still image. In this paper, we are interested in *inferring* the temporal dynamics from a previously unseen face image, and in the use of this information to further enhance the performance of Action Unit (AU) intensity estimation and dimensional affect estimation networks. To this end, we propose an image-to-image translation approach where a network is trained to generate a Dynamic Representation (DR) from a given image. Inspired by the dynamic image [1], our representation is formulated as a kernel that, when projected onto the adjacent frames, can sort them in time. It is therefore designed to be a single three-channel spatial data structure (similar to an RGB image), tasked with summarizing the motion that surrounds a given frame. This representation can be directly used in CNN-based networks for the tasks of AU intensity estimation and dimensional affect estimation.

During training, we are given a set of preceding and proceeding frames for each face image (i.e. sequences), from which the temporal evolution of adjacent frames can be learned in a self-supervised manner without using target representations (see Fig. 2). The network is then trained to generate a representation that, when projected onto each adjacent frame within a given window, is capable of sorting them in time. In other words, *we do not compute a set of target representations to learn our network*. In addition, we note that the temporal symmetry of facial actions could yield ambiguities when sorting frames in a strictly ascending order, and propose to sort frames relative to their distance to the central frame. We first validate empirically that the learned representation does have the capacity to sort adjacent frames in the temporal domain, thus illustrating its ability to capture short-term dependencies. Then, we show that our proposed approach can generate representations that are highly suitable for both AU intensity estimation and dimensional affect estimation. In particular, we propose a multi-level dynamic representation approach (see Fig. 1), that combines DRs generated for different temporal scales. We show this approach suffices to attain state of the art results in both tasks.

Our contributions can be summarized as follows:

- We propose an image-to-image translation network, tasked with *inferring* a dynamic representation of a given still image, designed to summarize the short-term motion surrounding it.
- We propose to train the network with a Rank Loss, enforcing the generated representations to rank both past and future frames according to their relative distance to the central frame. This way, the network not only learns to map an image to a corresponding representation, but also contributes to define it.
- We show how the inferred representations effectively summarize motion, and show how their use in combination with a given frame reaches state of the art results in the tasks of facial Action Unit (AU) intensity estimation and dimensional affect estimation.

2. Related Work

This Section reviews the closely related work, which we define as works related to the temporal modeling of facial motion and motion estimation from still images (our main goal), as well as image-to-image translation, image-based dynamic representations and self-supervised learning.

Temporal modeling of facial expressions Exploiting the temporal modeling of facial expressions on video sequences is a longstanding problem in Computer Vision. Some works have proposed to summarize short-term motion at the feature level, extending hand-crafted features to what is known as Three Orthogonal Planes (TOP) [4, 18]. Other works have exploited the use of a Fourier Transform [41], or spatio-temporal convolution [16, 51]. The majority of related work focus on using recurrent or latent-based models, in particular Recurrent Neural Networks (RNNs [6, 16, 51, 26, 10, 21]).

Motion Estimation Our work is related to motion prediction, where the goal is to infer motion from either still images or sequences. In this sense, the goal is to predict *what is going to happen next*. Some works have tackled this problem by predicting optical flow from still images [30, 46]. Others have proposed to infer the next frame to follow a preceding video sequence [3, 49]. In particular, [33] proposed to infer the next dynamic image, as it better correlates with the preceding frames. These methods do not attempt to summarize motion, but rather predict the most likely frame to follow a given image or image sequence.

Image-to-image translation Our work can be viewed as image-to-image translation, where a dynamic representation (a 3-channel image in our case) is generated from an input image. Works in image-to-image translation generally attempt to modify an input image to generate an output according to a target attribute or style, and thus do not have as a goal to *infer* any information *from* the input image [5, 15, 23, 47, 50, 53, 31]. These approaches generally rely on the use of Generative Adversarial Networks (GANs) [13], or any of its extensions [5, 15, 50]. GANs are a powerful tool to capture the target distribution, enforcing the networks to produce plausible outputs. However, as we

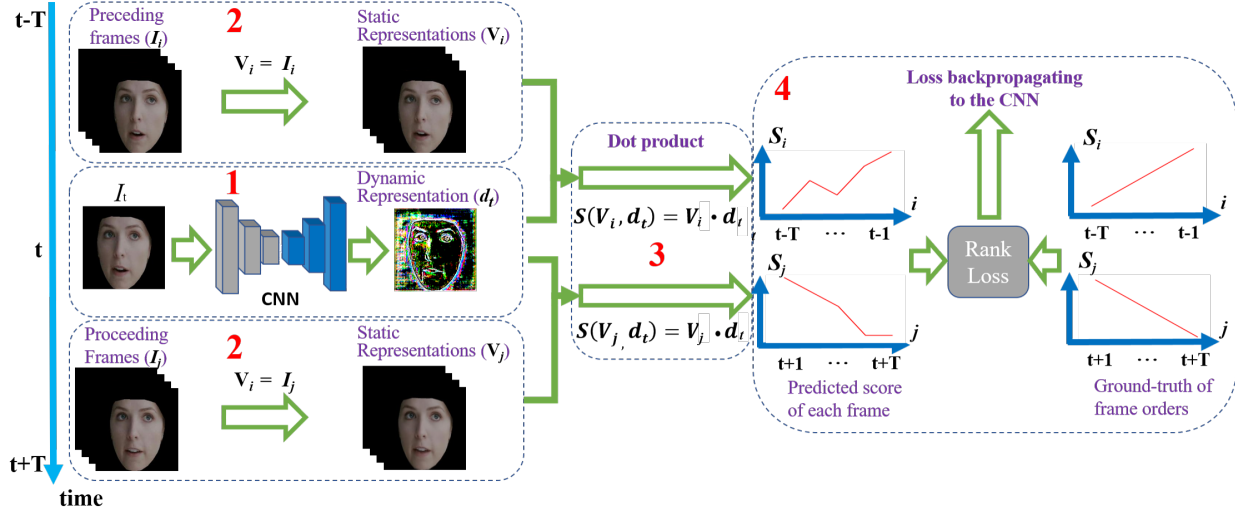


Figure 2. During training, we are given a set of sequences from which we can have access to the adjacent frames of a given image I_t . In (1), the given image I_t is forwarded to the network we aim to learn, that produces a DR d_t . We can measure the ranking capabilities of d_t by projecting it onto the preceding and proceeding frames (2). To rank the frames, we compute the difference between pair-wise scores, each computed as a dot product between the generated DR and the corresponding preceding or proceeding frame (3). These scores are used to compute a Rank Loss, which allow us to measure the extent of which the current d_t is correctly ranking the frames within the sequence. We can backpropagate the Rank Loss w.r.t. the parameters of the network that has produced the DR d_t (4). This way, the network not only learns to produce a correct representation d_t , but also contributes to define it.

shall see, we will not be using explicit target representations to learn our network, and therefore the use of GANs is not a suitable tool to learn the dynamic representations.

Dynamic Representation The basis of our work is referred to as dynamic representation. A dynamic representation is built so that it can rank order frames according to their position in a temporal sequence. The reasoning behind such an abstract representation is that if a representation has the power to rank all frames according to their temporal position in an image sequence, then it is a good descriptor of it, and thus can be used for machine learning tasks that require this temporal information. This hypothesis was validated for the task of human action recognition [1, 12]. The dynamic representation (referred to as *dynamic image* in [1]) was first presented as a short-term feature descriptor of image sequences. To obtain a dynamic image, one needs to *learn* it at test time from the set of frames that make up a sequence using e.g. RankSVM [40]. The use of RankSVM was further extended and converted into a pooling layer [12]. As we shall see, our network will be able to generate a dynamic representation *from still images*, that effectively summarizes not only past frames, but also future frames. Finally, it is worth mentioning that before the dynamic image, other methods were proposed to learn dynamic representations from image sequences, such as optical flow or Motion History Image [2].

Self-supervised Learning In this paper we propose to learn without explicitly generating target representations.

Instead, we will make use of a proxy loss function, called a Rank Loss, to train our network in a self-supervised manner. Self-supervised learning avoids the need of explicit target data, and instead explores the structure of the training data to supervise the training process, using e.g. temporal relations or semantic structures [8]. Some works on self-supervised learning have already used the temporal order of video frames to train networks, aiming to learn video representations of asymmetric human actions [11, 25] or analyze temporal coherence [17, 14]. To the best of our knowledge, we are the first to propose the use of a Rank Loss function to learn a dynamic representation of facial expressions in a self-supervised manner.

3. Proposed approach

Our goal is to train a network that generates a simple dynamic representation (DR) from a single face image, summarizing the motion around it. This is possible because facial actions are constrained by anatomy and behaviour causes strong correlations between adjacent frames. In this Section, we first define the goal of the DRs. To do so, we introduce the challenges that facial expressions pose for the task of learning a representation that summarizes motion, and propose an alternative representation that overcomes them. Then, we show how this representation can be learned. We note that generating a set of target representations to be the basis of a one-to-one mapping, common in image-to-image translation methods, is suboptimal, and

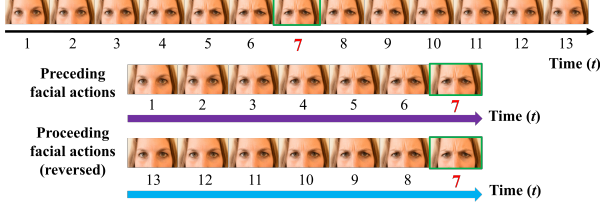


Figure 3. Temporal symmetry of facial motions. In the top row, a sequence of frames displaying a facial expression. With regards to the central frame $t = 7$, we can see that the preceding frames would look alike the proceeding frames, when the latter are reverted on time, i.e. the temporal pattern evolved from $t = 1$ to $t = 7$ with the natural arrow of time would be similar to the temporal pattern evolved from $t = 13$ to $t = 7$, under a reverted arrow of time.

propose an alternative approach to allow the learning process to be self-supervised by the temporal correlation between frames in the training set. Finally, we show how to apply the trained models for the task of AU intensity estimation and dimensional affect estimation. The full learning process is shown in Fig. 2.

3.1. Dynamic Representations

As mentioned above, our goal is to train a network that can infer a DR from a still face image. The first step consists of defining what output the network is expected to return for a given image. To do so, we draw our attention to the dynamic image algorithm [1]. The dynamic algorithm was presented as a kernel targeted with summarizing sequences of frames, with the goal of inferring the action on them. The main motivation behind generating a kernel that can rank frames relies on the assumption that such a representation should be a good descriptor of the sequence, as it encodes the temporal evolution of it. This assumption was empirically validated in [1]. In our framework, we are instead interested in having a per-frame representation, and therefore the original dynamic algorithm seems not to be an adequate representation. We empirically validate in Sec. 4.3 that, when the dynamic image is used as a basis to learn the representation, along with a reconstruction loss, the performance degrades substantially. In addition, it is important to note that facial expressions display a **symmetric** temporal pattern. As shown in Fig. 3, facial expressions are in many cases indistinguishable from their (temporally) reverted counterparts. In such cases, using the dynamic algorithm would incur in ambiguous representations for a given frame, according to whether it belongs to the “activation” part of an expression, or whether it belongs to its “deactivation” part. In order to overcome such limitation, we propose a DR that is targeted with ranking not only the preceding frames, but also the proceeding frames. In other words, the DR is chosen to be a kernel that can rank both past and future frames, based

on their temporal positions relative to the given face image. This way, the modeling of symmetric patterns is addressed in a more efficient way.

Let $I_t \in \mathbb{R}^{m \times n}$ be the given face image, and let $I_{t-T}, I_{t-T+1}, \dots, I_{t-1}$ and $I_{t+1}, I_{t+2}, \dots, I_{t+T}$ be the frames corresponding to a window of $2T + 1$ frames, centered at I_t . Let $V_a, a \in [t - T, t + T]$ be the static representation for the frame a . While in [1] V_a was defined as the cumulative feature representation of the frame a (i.e. $V_a = \sum_{a' < a} \phi(I_{a'})$, with ϕ a feature representation), in this paper we directly choose V_a to be the image itself, i.e. $V_a = I_a$, to avoid the feature representation to depend on other frames. Our goal is to generate a DR for frame I_t , namely d_t , with size equal to the static representation V_a , that can rank preceding and proceeding frames based on their relative temporal distance to I_t . The ranking of frames is performed by assigning a score to each, which is defined as the (Frobenius) inner product between the DR d_t and the representation V_a , with $a \in [t - T, t + T]$. Mathematically speaking, the score for frame $a \in [t - T, t + T]$, assigned by the DR d_t is defined as $S(d_t, V_a) = \langle d_t, V_a \rangle$. The scores are then used in an ordinal manner to sort the frames within the given window. In particular, we are interested in assigning ascending scores for frames $a < t$, and descending scores for $a > t$. This way, we can define the difference between the scores computed at time a and b as

$$\delta_{ab}(t) \doteq S(d_t, V_a) - S(d_t, V_b) \quad (1)$$

where a is chosen to be closer to t than b , with both a and b being either preceding or proceeding frames w.r.t. t . Then, our goal is to learn a DR d_t that makes $\delta_{ab}(t) > 0$ for pairs a, b corresponding to $a, b > t$ or $a, b < t$, with $|a - t| < |b - t|$. As we shall see in Sec. 3.2, we will only consider the actual value of $\delta_{ab}(t)$ when the pair a, b has been incorrectly ranked, i.e. when $\delta_{ab}(t) < 0$. Thus, we can define our target representation as a d_t that meets the following criteria:

$$\delta_{ab}(t) > 0 \quad \text{for} \quad \begin{cases} |a - t| < |b - t| \\ (a - t)(b - t) > 0 \end{cases} \quad (2)$$

It is important to remark that we compute the scores for the cases in which both frames are either before the current frame t , or after it. We are interested in computing ascending scores when $a, b < t$, and descending scores when $a, b > t$, and thus the cases where e.g. $a > t$ and $b < t$ would raise complex definitions.

3.2. Learning DR with Rank Loss

We now focus on how to learn the DR described above. Recall that our ultimate goal is to learn a static-to-dynamic projection f from a static face image I_t to its DR d_t , i.e. $d_t = f(I_t)$. This representation is tasked with meeting

the aforementioned criteria for all the training set of available frames. A priori, this could be accomplished by first generating the target representations from a sequence and then by training a network from the corresponding pairs, i.e. the centre image of the sequence and the target representations. In such case, one could use a reconstruction loss between the generated output and the corresponding target, so that the network learns to reproduce such a representation. However, we observe that when using a pre-defined representation to train the network, the generated outputs lack of generalization. We validate this empirically in Sec. 4.3. In other words, the network is forced to minimize a reconstruction loss w.r.t. a fixed representation, and therefore it does not take into account the capabilities of the generated output at the task of ranking adjacent frames. That is to say, *subtle errors in the reconstruction loss do not necessarily correlate with errors in the ranking of frames*. Instead, we want the network to also help design the DR.

In particular, instead of generating target representations, we propose to learn the DR by enforcing the network to produce outputs that directly meet Eqn. 2. More specifically, when the network generates an output for a given image, we project it onto the preceding and proceeding frames within a window of $N = 2T + 1$ frames, and compute the pair-wise scores using Eqn. 1 and Eqn. 2. Then, in a similar fashion to that of the RankSVM algorithm, we only account for the error committed by the pairs that have been incorrectly ranked. In addition, we add a rank success factor θ , to avoid small errors to be considered in the total loss. Mathematically speaking, let I_t be the given frame, corresponding to the central image of a window of $N = 2T + 1$ frames. Let $d_t = f(I_t)$ be the output of the network for the given frame. We want the generated DR to minimize the following rank loss function:

$$\begin{aligned}
L_f(d_t) = & \gamma \times \|d_t\|^2 - \varepsilon \\
& + \sum_{b=t-T}^{t-1} \sum_{a=b+1}^t \max(0, \theta - \delta_{ab}(t)) \\
& + \sum_{a=t}^{t+T-1} \sum_{b=a+1}^T \max(0, \theta - \delta_{ab}(t))
\end{aligned} \quad (3)$$

where recall $\delta_{ab}(t) = S(d_t, V_a) - S(d_t, V_b)$. In Eqn. 3, γ is a regularization factor, and ε is used as a relaxation factor to set an upper bound to the loss to avoid it to return extremely large values. The loss $L_f(d_t)$ can be differentiated w.r.t. the parameters of the network f , and therefore the network can be learned through typical backpropagation methods. The training process is also illustrated in Fig. 2. We want to recall that by using the Rank Loss function of Eqn. 3, we are not backpropagating w.r.t. a defined ‘‘ground-truth’’ d_t^* . In other words, our method is trained without the need of explicitly generating a target DR for each training image.

This way, the network also contributes to define the form of the DR.

3.3. Face analysis using Multi-level DR

We now describe how the DR shown above can be applied to face-related tasks. In particular, we observe that we can generate a multi-level set of DRs, each capturing a different temporal scale by using a different window length. We will validate that this combination allows to reach state of the art results in the tasks of AU intensity estimation and dimensional affect estimation. We first note that the generated d_t are 3-channel tensors, no matter the choice of T . Thus, we can train a different model for different values of T , and combine the outputs before applying them to further related tasks. Herein, we will explore the use of a **Single Dynamic Representation (SDR)**, using just the generated DR, and the use of a **Multi-level Dynamic Representation (MDR)**, which combines the output of networks trained using different time lengths T .

While there is no limit as to how many different levels can be used for further tasks, we want to keep the input network as simple as possible. To this end, in this paper we explore two different configurations, one for the AU intensity estimation and one for the dimensional affect estimation. We leave for future work to explore possible configurations. We rely on expert knowledge for each of the tasks: we choose a two-level representation for the AU intensity estimation task, and a three-level representation for the dimensional affect estimation task. In both, the first level corresponds to $T = 0$, i.e. a window length of one frame. In practice, this level does not require the training of a DR, as it basically consists of the input image. Indeed, we experimentally validate that the use of the input image along with the DR helps capturing the rich appearance details in the input image along with the temporal dynamics given by the DR. For the dimensional affect estimation task, we choose $T = 3$ for the second level. Finally, for both the second level of the AU intensity task and for the third level of the dimensional affect estimation task, we choose $T = 5$. Note that the total window size is defined as $N = 2T + 1$, i.e. $N = 7$ frames for $T = 3$, and $N = 11$ frames for $T = 5$. In all cases, we use a stride of $S = 2$ frames, i.e. we use every other frame to train the corresponding DR network. The chosen window then spans a set of $2TS + 1$ frames. This MDR yields a 6-channel tensor for the AU intensity task, and a 9-channel tensor for the dimensional affect estimation task. Fig. 1 shows a description of the MDR for this three level approach.

4. Experiments

To validate the proposed approach, we first evaluate the ranking capability of the generated DRs. Then, we demonstrate their value for the tasks of AU intensity and dimen-

sional affect estimation. We will show that our MDR produces state of the art results in both tasks. It is important to note that, at test time, a single face image is used as the input to generate the corresponding MDR, and this representation is used to predict the values of the corresponding task. In addition, we want to remark that the MDR network is trained using the RECOLA dataset (see below), whereas the AU intensity network and the dimensional affect estimation network are trained using the BP4D and SEMAINE datasets, respectively. In other words, the DR networks are trained using a different dataset to those used to the corresponding tasks.

4.1. Datasets

Experiments were conducted on three face datasets: RECOLA [32], SEMAINE [24] and BP4D [54]. *The RECOLA dataset is solely used to train the DR networks, whereas SEMAINE and BP4D are used to test both the capabilities of the DR to rank the corresponding frames, and to train and test the corresponding face-related tasks.* For the RECOLA dataset, we use the 27 videos corresponding to the AVEC 2016 challenge [42], each containing approximately 5 minutes of people performing video conference. For the SEMAINE dataset, we use the subset predefined by the AVEC 2012 challenge [38], which consists of 31 training videos, 32 validation videos, and 32 test videos. All frames have been annotated with valence and arousal intensities, each lying in the range $[-1, 1]$. For the BP4D dataset, we use the partitions predefined by the FERA 2015 challenge [43]. There are 75,586 frames for training, 71,260 frames for development and 75,726 frames in the test set. All the frames have been annotated for five AUs (AU6, AU10, AU12, AU14, and AU17), each lying in the range $[0, 5]$. The sampling rate of RECOLA and BP4D was 25 fps, whereas the sampling rate of SEMAINE was ~ 50 fps. For this reason, wherever we refer to the stride S when setting up the span of frames to be considered, this will be automatically scaled to $2S$ for the SEMAINE dataset.

4.2. Implementation details

Configuration All the experiments are carried out using the PyTorch library [29] for deep learning. The network chosen for the task of generating the DRs is the UNet [34]. Both the input and the output of the UNet are tensors of size $224 \times 224 \times 3$. The parameter θ in Eqn. 3 is set beforehand to ensure the chance level ranking accuracy to be less than 0.1%. For the task of AU recognition, we re-trained the network proposed by [36], as multi-task learning has been frequently adopted for AU recognition [22, 27]. For affect estimation, we followed the setting in [21], and fine-tuned the VGG-16 face network [28], whose last layer is modified to have an output size of 1.

Pre-processing We used the publicly available iCCR

face tracker of [35] to first detect a set of 66 facial landmarks. Using these landmarks, images are cropped to meet the network size. Then, all pixels corresponding to the outer part of the convex hull defined by the landmarks are set to zero to remove all non facial appearance information.

Training details The UNet was trained using an Adam optimizer [7] with a learning rate of 10^{-3} , and $\beta = (0.5, 0.9)$. During the ranking experiment, the UNet is trained and validated using the RECOLA dataset, and tested on all frames in the BP4D and SEMAINE datasets. For the face related tasks, we utilized the trained DR models to generate the DRs for each frame in the SEMAINE and BP4D datasets, and then trained the corresponding AU/dimensional affect models using the generated representations.

4.3. Frame Ranking

In this section, we report the ability of the generated DRs to rank the corresponding adjacent frames in SEMAINE and BP4D. In particular, we evaluate the capabilities of our model under different scenarios by choosing a set of different window lengths and strides to train the networks and generate the corresponding DRs. The number of frames used per training image is $N = 2T + 1$ (T preceding frames, T proceeding frames and the given frame). We sample N frames using four different strides S . The image sequence range is then of $N \times S$ frames. We check the capabilities of our network for $T = \{3, 5, 7, 9\}$ (i.e. $N = \{7, 11, 15, 19\}$). In the most extreme case, i.e. when $T = 9$ and $S = 4$, the ranking is measured on a window size of $N = 2T + 1 = 19$ frames, evenly sampled from a sequence of $N \times S = 76$ frames. At test time, frames are chosen following the same sampling procedure as that of the corresponding model. To compute the ranking accuracy, we compute the DR for each of the images available in the corresponding datasets. Then, we project the generated DR onto the frames lying within the corresponding window of N frames, sampled with a stride S . We then measure the distances $\delta_{ab}(t)$ as defined in Eqn. 1, and measure the percentage of pairs that are correctly ranked, i.e. the percentage of pairs for which $\delta_{ab}(t) > 0$.

In Fig. 4 we report the ranking accuracy of our method, measured as a percentage of correctly ranked frames w.r.t. total number of frames evaluated. We report the average accuracy measured across both SEMAINE and BP4D. The results shown with a dash line correspond to applying a RankSVM at test time, trained using the frames that were later ranked by it. The results given by the RankSVM are treated as an *upper bound* for the ranking accuracy.

We compare the accuracy of our method against methods trained using target representations obtained by the Dynamic Image algorithm, which generates an explicit target representation for each training image using the RankSVM

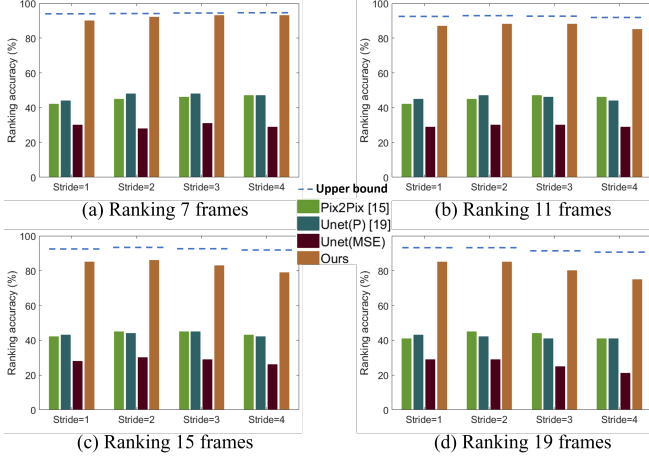


Figure 4. Average ranking accuracy (%) on two datasets. Four generative models are trained using RECOLA dataset and tested on SEMAINE and BP4D datasets. RankSVM classifiers were trained on SEMAINE and BP4D datasets, and each classifier only rank its training frames. The results obtained by RankSVM are treated as the upper bound.

algorithm [1]. We generated the corresponding dynamic images both forward (**DI**), and backward (**BDI**) in time, and applied them as the targets to train several generative networks. We trained a different network for the DI and the BDI images, respectively. In particular, Fig. 4 shows the results given by the following approaches:

- UNet (MSE). Using the dynamic images, we train the model using as objective the Mean Squared Error.
- UNet (P). In this method, we use the dynamic images as target representation, and the objective function used to train the model is the Perceptual loss proposed in [19].
- Pix2Pix [15] refers to using a conditional GAN, again using the dynamic images as the corresponding targets.

The results shown in Fig. 4 show how our method achieves similar results to those given by the RankSVM, which uses at test time the adjacent frames to compute the kernel. Remarkably, our method yields around 80% accuracy even for the longest cases (i.e. when ranking 19 frames with different strides). In addition, we can see that, when pairing the input images with a DR to serve as a basis to learn our network, the ranking accuracy at test time degrades substantially. This illustrates the contribution of the Rank Loss at the task of defining the form of the DR, allowing for a better generalization. An example is shown in Fig. 5. It can be seen that our proposed DR is capable of accurately ranking both preceding and proceeding frames, something not possible with methods trained with explicit target representations. In addition, the RankSVM would

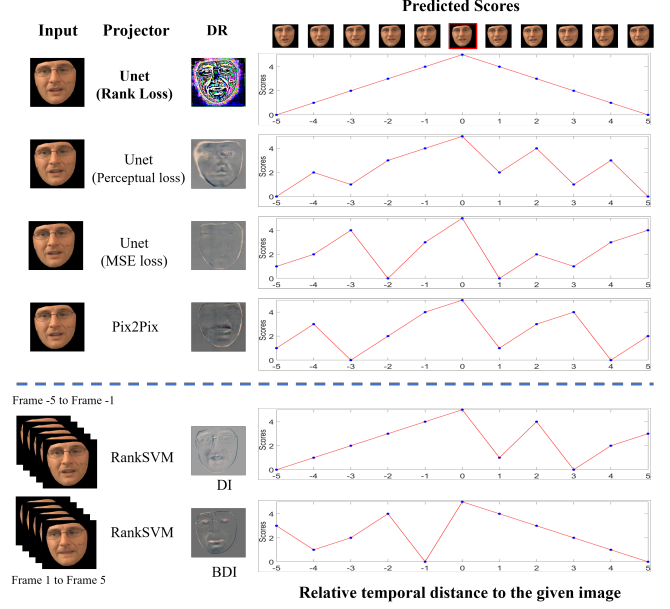


Figure 5. Examples of ranking frames using DRs generated by different methods. The networks of Pix2Pix, Unet(P) and Unet(MSE) were trained using Seq DI as the target.

only be able to rank preceding or proceeding frames, and it would need to be given the temporal information at test time. This clearly illustrates the capacities of our proposed approach at the task of sorting the surrounding frames of a given image relative to their position to it.

4.4. Face-related tasks

In this section, we show the efficacy of the proposed approach for the tasks of AU intensity estimation and dimensional affect estimation, for which the temporal modeling would generally enhance the performance. As we demonstrate, our still image-based method yields similar results to those given by the models using image sequence, which make use of the temporal information. However, contrary to the temporal modeling of adjacent frames, our method only requires a single still image.

AU intensity estimation We first evaluate the contribution of our proposed approach at the task of AU intensity estimation. To do so, we re-train the network proposed by [36] using the training partition of BP4D. Then, we evaluate the performance of the proposed SDR and MDR approaches against using a single input image, as in [36]. We evaluate each method employing the same measures as those used by the FERA 2015 [43] to rank participants: the Intra-Class Correlation (ICC(3,1), [39]), and the Mean Squared Error (MSE). We also compare the performance of our method against using the alternatives shown in Sec. 4.3, i.e. those using an explicit target DR for training. Finally, we com-

	AU	6	10	12	14	17	Avg.
ICC	CCNN-IT [44]	0.75	0.69	0.86	0.40	0.45	0.63
	2DC [22]	0.76	0.71	0.85	0.45	0.53	0.66
	VGP-AE [9]	0.75	0.66	0.88	0.47	0.49	0.65
	HG-HMR [36]	0.79	0.80	0.86	0.54	0.43	0.68
	Pix2Pix* [15]	0.59	0.62	0.68	0.29	0.31	0.50
	UNet(P)* [19]	0.55	0.65	0.65	0.30	0.26	0.48
	UNet(MSE)*	0.56	0.63	0.66	0.29	0.26	0.48
	SDR+HG-HMR	0.78	0.80	0.85	0.47	0.45	0.67
	MDR+HG-HMR	0.77	0.83	0.87	0.62	0.49	0.72
MSE	CCNN-IT [44]	1.23	1.69	0.98	2.72	1.17	1.57
	2DC [22]	0.75	1.02	0.66	1.44	0.88	0.95
	VGP-AE [9]	0.82	1.28	0.70	1.43	0.77	1.00
	HG-HMR* [36]	0.77	0.92	0.65	1.57	0.77	0.94
	Pix2Pix* [15]	1.22	1.31	0.85	1.90	0.92	1.24
	UNet(P)* [19]	1.53	1.08	1.07	1.62	0.95	1.25
	UNet(MSE)*	1.09	1.55	1.18	2.12	1.15	1.42
	SDR+HG-HMR	0.88	0.84	0.75	1.90	0.60	0.99
	MDR+HG-HMR	0.99	0.79	0.64	1.34	0.48	0.85

Table 1. AU intensities estimation results on BP4D dataset. * denotes results obtained by our own implementation

pare our approach against most recent works reporting state of the art results on the BP4D dataset. The results are shown in Table 1.

From the results shown in Table 1, it can be seen that our MDR approach is capable of achieving state of the art results even when using as input a single image. In addition, we observe that the SDR achieved similar results to using a single image, i.e. the SDR approach gives similar results to those of [36]. We conjecture that while the SDR can encode the temporal pattern around a given frame, the original input image is rich in appearance details that remain important to infer the AU intensities. Thus, the best results are attained when combining the input image with the DRs (i.e. using the MDR approach). In other words, the estimated dynamics helped the still image-based AU intensity estimation to yield better results.

Dimensional affect estimation We measure the performance of our proposed approach at the task of dimensional affect estimation, i.e. at predicting the values of valence and arousal. To do so, we use the standard measures reported on the SEMAINE dataset, i.e. the Pearson Correlation Coefficient (PCC), and the Mean Squared Error (MSE). As introduced in Sec. 4.2, we fine-tune the VGG-16 network [28] for each of the alternatives aforementioned. More specifically, we use a VGG-16 network for the inputs generated by our SDR and MDR, as well as for the methods trained using fixed representations, UNet(MSE), UNet(P), and Pix2Pix. In addition, we compare our approach against using a single image as input to the VGG-16 (SI+VGG). The results are shown in Table 2. Again, we compare against state of

	Arousal		Valence	
Method	PCC	MSE	PCC	MSE
Savran et al. [†] [37]	0.251	N.A.	0.210	N.A.
Kaltwang et al. [†] [20]	0.310	0.042	0.310	0.058
Zhang et al. [52]	0.070	N.A.	0.241	N.A.
SI+VGG * [28]	0.246	0.056	0.258	0.084
Pix2Pix*+VGG [15]	0.091	0.166	0.088	0.195
UNet(P)*+VGG [19]	0.192	0.125	0.134	0.172
UNet(MSE)*+VGG	0.063	0.181	0.082	0.189
SDR+VGG	0.306	0.078	0.299	0.082
MDR+VGG	0.335	0.058	0.316	0.072

Table 2. Affect estimation results on the SEMAINE dataset. SI denotes the still face image; * denotes our own implementation; [†] denotes methods that rely on the use of temporal information at test time

the art methods on valence and arousal estimation using visual information, including those that make use of temporal information to improve their performance ([37, 20]). As it can be seen, our proposed MDR approach attains state of the art results using as input only still images.

5. Conclusion

In this paper, we have proposed a novel approach to model temporal dynamics from static face images. This approach allows us to train a network to infer a DR for a previously unseen test image, which effectively summarize dynamics surround it. We illustrated that the generated DRs can be used indistinctly for the tasks of Action Unit intensity and dimensional affect estimation, attaining state of the art results. We empirically validated the capacity of the DRs to rank unseen frames in test time, as well as their contribution to the face-related tasks. In addition, we validated that a network trained with a rank loss function generalizes better to unseen images than a model trained using pre-defined representations, i.e. we demonstrated the ability of our network to be properly learned without the need of target representations. Experimental results have shown that our method is powerful not only for the task of ranking adjacent frames in different facial actions but also for the tasks of estimating the AU and affect intensities from the face.

References

- [1] H. Bilen, B. Fernando, E. Gavves, and A. Vedaldi. Action recognition with dynamic image networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.
- [2] A. F. Bobick and J. W. Davis. The recognition of human movement using temporal templates. *IEEE Transactions on pattern analysis and machine intelligence*, 23(3):257–267, 2001.

- [3] B. Chen, W. Wang, and J. Wang. Video imagination from a single image with transformation generation. In *Proceedings of the on Thematic Workshops of ACM Multimedia 2017*, pages 358–366. ACM, 2017.
- [4] J. Chen, Z. Chen, Z. Chi, and H. Fu. Facial expression recognition in video with multiple feature fusion. *IEEE Transactions on Affective Computing*, 9(1):38–50, 2018.
- [5] Y. Choi, M. Choi, M. Kim, J.-W. Ha, S. Kim, and J. Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. *arXiv preprint*, 2018.
- [6] W.-S. Chu, F. De la Torre, and J. F. Cohn. Learning spatial and temporal cues for multi-label facial action unit detection. In *2017 12th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2017)*, pages 25–32. IEEE, 2017.
- [7] J. B. Diederik P. Kingma. Adam: A method for stochastic optimization. In *Int’l Conference for Learning Representations (ICLR)*, 2015.
- [8] C. Doersch, A. Gupta, and A. A. Efros. Unsupervised visual representation learning by context prediction. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1422–1430, 2015.
- [9] S. Eleftheriadis, O. Rudovic, M. P. Deisenroth, and M. Pantic. Variational gaussian process auto-encoder for ordinal prediction of facial action units. In *Asian Conference on Computer Vision*, pages 154–170. Springer, 2016.
- [10] Y. Fan, X. Lu, D. Li, and Y. Liu. Video-based emotion recognition using cnn-rnn and c3d hybrid networks. In *Proceedings of the 18th ACM International Conference on Multimodal Interaction*, pages 445–450. ACM, 2016.
- [11] B. Fernando, H. Bilen, E. Gavves, and S. Gould. Self-supervised video representation learning with odd-one-out networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3636–3645, 2017.
- [12] B. Fernando, E. Gavves, J. Oramas, A. Ghodrati, and T. Tuytelaars. Rank pooling for action recognition. *IEEE transactions on pattern analysis and machine intelligence*, 39(4):773–787, 2017.
- [13] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [14] R. Goroshin, J. Bruna, J. Tompson, D. Eigen, and Y. LeCun. Unsupervised learning of spatiotemporally coherent metrics. In *Proceedings of the IEEE international conference on computer vision*, pages 4086–4093, 2015.
- [15] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5967–5976, 2017.
- [16] S. Jaiswal and M. Valstar. Deep learning the dynamic appearance and shape of facial action units. In *Applications of Computer Vision (WACV), 2016 IEEE Winter Conference on*, pages 1–8. IEEE, 2016.
- [17] D. Jayaraman and K. Grauman. Slow and steady feature analysis: higher order temporal coherence in video. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3852–3861, 2016.
- [18] B. Jiang, M. F. Valstar, B. Martinez, and M. Pantic. A dynamic appearance descriptor approach to facial actions temporal modeling. *IEEE Trans. Cybernetics*, 44(2):161–174, 2014.
- [19] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision*, pages 694–711. Springer, 2016.
- [20] S. Kaltwang, S. Todorovic, and M. Pantic. Doubly sparse relevance vector machine for continuous facial behavior estimation. *IEEE transactions on pattern analysis and machine intelligence*, 38(9):1748–1761, 2016.
- [21] D. Kollias, P. Tzirakis, M. A. Nicolaou, A. Papaioannou, G. Zhao, B. Schuller, I. Kotsia, and S. Zafeiriou. Deep affect prediction in-the-wild: Aff-wild database and challenge, deep architectures, and beyond. *International Journal of Computer Vision*, Feb 2019.
- [22] D. Linh Tran, R. Walecki, S. Eleftheriadis, B. Schuller, M. Pantic, et al. Deepcoder: Semi-parametric variational autoencoders for automatic facial action coding. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3190–3199, 2017.
- [23] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015.
- [24] G. McKeown, M. Valstar, R. Cowie, M. Pantic, and M. Schroder. The semaine database: Annotated multimodal records of emotionally colored conversations between a person and a limited agent. *IEEE Transactions on Affective Computing*, 3(1):5–17, 2012.
- [25] I. Misra, C. L. Zitnick, and M. Hebert. Shuffle and learn: unsupervised learning using temporal order verification. In *European Conference on Computer Vision*, pages 527–544. Springer, 2016.
- [26] B. H. Mohammad Mahoor et al. Facial expression recognition using enhanced deep 3d convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 30–40, 2017.
- [27] J. Nicolle, K. Bailly, and M. Chetouani. Facial action unit intensity prediction via hard multi-task metric learning for kernel regression. In *2015 11th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG)*, volume 6, pages 1–6. IEEE, 2015.
- [28] O. M. Parkhi, A. Vedaldi, A. Zisserman, et al. Deep face recognition. In *BMVC*, volume 1, page 6, 2015.
- [29] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in pytorch. 2017.
- [30] S. L. Pintea, J. C. van Gemert, and A. W. M. Smeulders. Deja vu: Motion prediction in static images. In *ECCV*, 2014.
- [31] A. Pumarola, A. Agudo, A. M. Martinez, A. Sanfeliu, and F. Moreno-Noguer. Ganimation: Anatomically-aware facial animation from a single image. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 818–833, 2018.

- [32] F. Ringeval, A. Sonderegger, J. Sauer, and D. Lalanne. Introducing the recola multimodal corpus of remote collaborative and affective interactions. In *Automatic Face and Gesture Recognition (FG), 2013 10th IEEE International Conference and Workshops on*, pages 1–8. IEEE, 2013.
- [33] C. Rodriguez, B. Fernando, and H. Li. Action anticipation by predicting future dynamic images. *arXiv preprint arXiv:1808.00141*, 2018.
- [34] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.
- [35] E. Sánchez-Lozano, G. Tzimiropoulos, B. Martinez, F. De la Torre, and M. Valstar. A functional regression approach to facial landmark tracking. *IEEE transactions on pattern analysis and machine intelligence*, 40(9):2037–2050, 2018.
- [36] E. Sanchez-Lozano, G. Tzimiropoulos, and M. Valstar. Joint action unit localisation and intensity estimation through heatmap regression. *BMVC*, 2018.
- [37] A. Savran, H. Cao, M. Shah, A. Nenkova, and R. Verma. Combining video, audio and lexical indicators of affect in spontaneous conversation via particle filtering. In *Proceedings of the 14th ACM international conference on Multimodal interaction*, pages 485–492, 2012.
- [38] B. Schuller, M. Valster, F. Eyben, R. Cowie, and M. Pantic. Avec 2012: the continuous audio/visual emotion challenge. In *Proceedings of the 14th ACM international conference on Multimodal interaction*, pages 449–456. ACM, 2012.
- [39] P. E. ShROUT and J. L. Fleiss. Intraclass correlations: uses in assessing rater reliability. *Psychological bulletin*, 86(2):420, 1979.
- [40] A. J. Smola and B. Schölkopf. A tutorial on support vector regression. *Statistics and computing*, 14(3):199–222, 2004.
- [41] S. Song, L. Shen, and M. Valstar. Human behaviour-based automatic depression analysis using hand-crafted statistics and deep learned spectral features. In *Automatic Face & Gesture Recognition (FG 2018), 2018 13th IEEE International Conference on*, pages 158–165. IEEE, 2018.
- [42] M. Valstar, J. Gratch, B. Schuller, F. Ringeval, D. Lalanne, M. Torres Torres, S. Scherer, G. Stratou, R. Cowie, and M. Pantic. Avec 2016: Depression, mood, and emotion recognition workshop and challenge. In *Proceedings of the 6th International Workshop on Audio/Visual Emotion Challenge*, pages 3–10. ACM, 2016.
- [43] M. F. Valstar, T. Almaev, J. M. Girard, G. McKeown, M. Mehu, L. Yin, M. Pantic, and J. F. Cohn. Fera 2015-second facial expression recognition and analysis challenge. In *2015 11th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG)*, volume 6, pages 1–8. IEEE, 2015.
- [44] R. Walecki, V. Pavlovic, B. Schuller, M. Pantic, et al. Deep structured learning for facial action unit intensity estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3405–3414, 2017.
- [45] J. Walker, C. Doersch, A. Gupta, and M. Hebert. An uncertain future: Forecasting from static images using variational autoencoders. In *European Conference on Computer Vision*, pages 835–851. Springer, 2016.
- [46] J. Walker, A. Gupta, and M. Hebert. Dense optical flow prediction from a static image. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2443–2451, 2015.
- [47] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, A. Tao, J. Kautz, and B. Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, page 5, 2018.
- [48] F. Xu, J. Zhang, and J. Z. Wang. Microexpression identification and categorization using a facial dynamics map. *IEEE Transactions on Affective Computing*, 8(2):254–267, 2017.
- [49] T. Xue, J. Wu, K. Bouman, and W. Freeman. Visual dynamics: Stochastic future generation via layered cross convolutional networks. *IEEE transactions on pattern analysis and machine intelligence*, 2018.
- [50] Z. Yi, H. R. Zhang, P. Tan, and M. Gong. Dualgan: Unsupervised dual learning for image-to-image translation. In *ICCV*, pages 2868–2876, 2017.
- [51] K. Zhang, Y. Huang, Y. Du, and L. Wang. Facial expression recognition based on deep evolutionary spatial-temporal networks. *IEEE Transactions on Image Processing*, 26(9):4193–4203, 2017.
- [52] L. Zhang, D. Tjondronegoro, and V. Chandran. Representation of facial expression categories in continuous arousal-valence space: feature and correlation. *Image and Vision Computing*, 32(12):1067–1079, 2014.
- [53] R. Zhang, P. Isola, and A. A. Efros. Colorful image colorization. In *European Conference on Computer Vision*, pages 649–666. Springer, 2016.
- [54] X. Zhang, L. Yin, J. F. Cohn, S. Canavan, M. Reale, A. Horowitz, P. Liu, and J. M. Girard. Bp4d-spontaneous: a high-resolution spontaneous 3d dynamic facial expression database. *Image and Vision Computing*, 32(10):692–706, 2014.