



# Self-supervised pain intensity estimation from facial videos via statistical spatiotemporal distillation

Mohammad Tavakolian<sup>a,\*\*</sup>, Miguel Bordallo Lopez<sup>b</sup>, Li Liu<sup>a,c</sup>

<sup>a</sup>Center for Machine Vision and Signal Analysis (CMVS), University of Oulu, Finland

<sup>b</sup>VTT Technical Research Centre of Finland Ltd, Oulu, Finland

<sup>c</sup>National University of Defense Technology, China

## ABSTRACT

Recently, automatic pain assessment technology, in particular automatically detecting pain from facial expressions, has been developed to improve the quality of pain management, and has attracted increasing attention. In this paper, we propose self-supervised learning for automatic yet efficient pain assessment, in order to reduce the cost of collecting large amount of labeled data. To achieve this, we introduce a novel similarity function to learn generalized representations using a Siamese network in the pretext task. The learned representations are finetuned in the downstream task of pain intensity estimation. To make the method computationally efficient, we propose Statistical Spatiotemporal Distillation (SSD) to encode the spatiotemporal variations underlying the facial video into a single RGB image, enabling the use of less complex 2D deep models for video representation. Experiments on two publicly available pain datasets and cross-dataset evaluation demonstrate promising results, showing the good generalization ability of the learned representations.

© 2020 Elsevier Ltd. All rights reserved.

## 1. Introduction

The concept of pain is learned by individuals through their life experiences. Hence, pain is a subjective phenomenon, which is influenced to varying degrees by biological, psychological, and social factors. Formally, pain can be described as an unpleasant sensory and emotional experience typically caused by actual or potential tissue injury. Therefore, pain is considered as an indicator for the health condition, which identifies harmful conditions for the body. Reliable understanding of pain contributes to disease diagnosis and provides useful information for healthcare personnel to choose adequate treatment to avoid long-lasting consequences for individuals' health.

Automatic pain assessment has attracted increasing attention of the research community, as it offers continue and objective pain assessment that improves the clinical outcome. Using visual data, many automatic pain assessment methods have been proposed by analyzing facial expressions which are a reliable

indicator of pain, estimating pain intensity, and distinguishing facial expressions caused by pain from other ordinary facial expressions. In automatic pain assessment, the development of powerful feature representations from facial videos plays a very important role and thus has been one main focus of research. Earlier works attempt to address this problem by either ignoring the temporal cue and extracting spatial features from frames, or treating videos as volumetric objects to capture spatiotemporal features. For instance, [Lucey et al. \(2012\)](#) used SVM classifiers to classify there levels of pain intensity. [Kaltwang et al. \(2012\)](#) computed a facial representation by extracting LBP and DCT features from video frames. [Zhao et al. \(2016\)](#) introduced an alternating direction technique of multipliers to solve Ordinal Support Vector Regression. Recent advances in deep learning, in particular Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), have brought significant progress for the problem of pain analysis from facial videos. State of the art results in this field are typically supervised deep learning methods. [Zhou et al. \(2016\)](#) trained a Recurrent Convolutional Neural Network (RCNN) to solve a regression problem. Similarly, [Rodriguez et al. \(2017\)](#) extracted features from frames using a CNN and

<sup>\*\*</sup>Corresponding author:

*e-mail:* [firstname.lastname@oulu.fi](mailto:firstname.lastname@oulu.fi) (Mohammad Tavakolian)

fed these features to a LSTM. Tavakolian and Hadid (2019) proposed a Spatiotemporal Convolutional Neural Network to capture low, mid, and long-range variations in the face for pain intensity estimation. However, data annotation in pain analysis is typically done using the Facial Action Coding System (FACS) (Ekman and Friesen (1978)) which firstly decomposes facial expressions into basic action units and then codes each action unit with onset, offset and intensity. FACS based labeling is costly and time consuming, since it takes around two hours for a FACS expert to code one minute video (Littlewort et al. (2009)).

In addition, the generalization ability of supervised learning is rather limited, *i.e.* defined largely by the training data. In other words, the performance of existing methods is limited to the data which they are trained on (Othman et al. (2019); Wang et al. (2018b); Yan (2016)). The differences in culture, gender and age, as well as imaging conditions such as lighting, pose, low resolution and noise are counted as factors for performance deterioration when generalizing methods across datasets (Dailey et al. (2010)). Othman et al. (2019) performed cross-dataset pain recognition using two datasets and their combination. However, they reduced the size of the deep model as well as training data by sampling frames that have high chance of pain occurrence. Currently, to the best of our knowledge, another challenging issue in facial video based pain analysis is the lack of large scale benchmark datasets due to the difficulties in data collection and labeling, while deep learning methods require large scale datasets for training to obtain good performance. Therefore, leveraging knowledge from pretrained models and using it to advance pain assessment could be a good option. In order to improve the generalization of pain assessment methods, we advocate the test of performance robustness by applying cross-database evaluation and the use of dataset combination to advance pain assessment.

To avoid time-consuming and costly data annotations, self-supervised learning has been put forward to learn visual representations from large-scale unlabeled data (Jing and Tian (2020)). It is concerned with learning semantically meaningful features from unlabeled data attributes. Hence, self-supervised learning can be considered as an alternative to supervised learning in automatic pain assessment, where collecting large-scale labeled dataset is a tedious task. In this paper, we propose a self-supervised learning method to capture visual features from unlabeled data using a pretext task to improve the generalization of learned representations for pain intensity estimation. To be specific, we learn feature representations based on measuring sample similarity of large-scale unlabeled data with a Siamese network (Bromley et al. (1994); Hu et al. (2014)). Then, we transfer the learned representations to the downstream task of estimate pain intensity level. We note that deep models usually have large computational complexity for spatiotemporal processing of videos. On the other hand, processing short intervals of videos is not efficient as they contain less information (Wang et al. (2018a)). To circumvent the deficiencies caused by processing videos on short intervals and/or avoiding deep models for temporal processing, we propose a highly simple, efficient yet effective method called

Statistical Spatiotemporal Distillation (SSD) to aggregate motion and appearance information underlying a video into one feature map. By temporally dividing the videos into smaller segments, SSD extracts the mutual information between consecutive segments to generate feature maps based on Gaussian Scale Mixture. Using the obtained feature maps (*i.e.* 2D representations of videos) allows employing less complex deep models, which are designed for images, for video processing.

In summary, our main contributions in this paper are summarized as follows.

- To the best of our knowledge, this is the first work that explores the idea of self-supervised representation learning in pain estimation. To achieve this, we adapt a novel similarity function through a Siamese network to learn representations from unlabeled data.
- In order to achieve computational efficiency, we propose a simple, efficient yet effective method (SSD) to encode statistical information of the underlying dynamic and appearance of an arbitrary-length video into a single image map.
- We extensively validate the effectiveness of our proposed method for pain intensity estimation in cross-dataset settings, using two publicly available datasets.

## 2. Proposed Method for Automatic Pain Assessment

In this section, we present firstly in detail the proposed SSD approach, which enables us to effectively aggregate the appearance and dynamic variations within a video into one single RGB image, and, the use of less complex 2D deep models to analyze the data. Then, we present how to explore self-supervised learning to learn generalized representations from unlabeled data for automatic pain assessment.

### 2.1. Statistical Spatiotemporal Distillation

Visual attention is usually given to the regions that have more descriptive information. Inspired by information theory, the local information of an image can be quantified in terms of sequences of bits (Hossein Khatoonabadi et al. (2015)). We extend this notion to the temporal dimension in order to capture the discriminative spatiotemporal information. To this end, under a Markovian assumption, we delineate a video as an image map by devising a statistical model for a neighboring group of pixels using Gaussian Scale Mixture (GSM).

Firstly, we divide a given video into several nonoverlapping segments. Then, the mutual information between consecutive segments (*i.e.* the total perceptual information content of consecutive segments) is explored to model the variations in dynamics and appearance. Hence, we encode the spatiotemporal variations of two consecutive segments into one image map. Further, we use an aggregation function to combine several image maps that are obtained by processing consecutive segments throughout the video (Figure 1).

Our approach is formally formulated as follows. A block volume,  $\mathbf{x}_i$ , in a video segment is modeled with GSM:  $\mathbf{x}_i =$

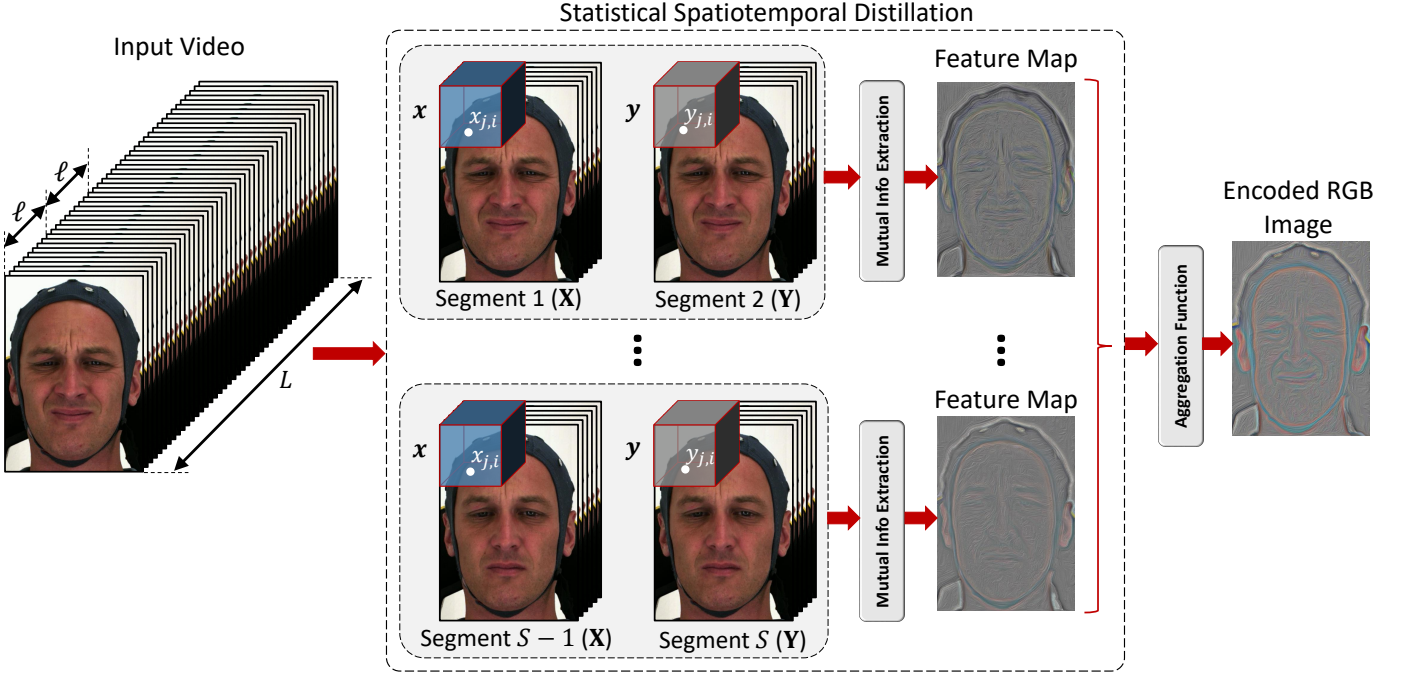


Fig. 1: The outline of our proposed SSD for video representation. Given a video of length  $L$ , we divide it into non-overlapping segments of smaller length  $\ell$ . SSD generates feature maps by computing mutual information between block volumes (spatiotemporal neighborhood) of consecutive segments. Using an aggregation function, we obtain an encoded RGB image from the input video.

$\alpha_i \mathbf{u}_i$ , where  $\mathbf{u}_i$  is a zero-mean Gaussian vector,  $\alpha_i$  is a mixing multiplier. The general form of GSM allows  $\alpha_i$  to be random variable that has a certain distribution in a continuous scale. For computation simplicity, we assume that  $\alpha_i$  only takes a fixed value at each block volume (*i.e.* it has different values in each block volume). For notation simplicity, we omit the subscript  $i$ . In practical scenarios, the noise effect needs to be taken into consideration in our model. Hence, we extend the model of the block volume as:

$$\mathbf{p} = \mathbf{x} + \mathbf{n}_1 = \alpha \mathbf{u} + \mathbf{n}_1, \quad (1)$$

where  $\mathbf{n}_1$  is Gaussian noise. Intuitively, each block volume  $\mathbf{p}$  undergoes spatiotemporal variations  $\mathbf{v}$  over time, leading to a distorted version  $\mathbf{q}$ :

$$\mathbf{q} = \mathbf{y} + \mathbf{n}_2 = g\alpha \mathbf{u} + \mathbf{v} + \mathbf{n}_2, \quad (2)$$

where  $\mathbf{y}$  represents deformation of  $\mathbf{x}$  (*i.e.*  $\mathbf{y} = g\alpha \mathbf{u} + \mathbf{v}$ ),  $g$  is a gain factor, and  $\mathbf{n}_2$  denotes Gaussian noise.  $\mathbf{n}_1$  and  $\mathbf{n}_2$  are independent with covariance matrices  $\mathbf{C}_{\mathbf{n}_1} = \mathbf{C}_{\mathbf{n}_2} = \sigma_n^2 \mathbf{I}$ , where  $\mathbf{I}$  is an identity matrix. Then, the covariance matrices of  $\mathbf{x}$ ,  $\mathbf{y}$ ,  $\mathbf{p}$ , and  $\mathbf{q}$  are derived straightforwardly as:

$$\mathbf{C}_x = \alpha^2 \mathbf{C}_u, \quad \mathbf{C}_y = g^2 \alpha^2 \mathbf{C}_u + \sigma_v^2 \mathbf{I} \quad (3)$$

$$\mathbf{C}_p = \alpha^2 \mathbf{C}_u + \sigma_n^2 \mathbf{I}, \quad \mathbf{C}_q = g^2 \alpha^2 \mathbf{C}_u + \sigma_v^2 \mathbf{I} + \sigma_n^2 \mathbf{I} \quad (4)$$

where  $\mathbf{C}_u$  is the covariance matrix of  $\mathbf{u}$ .

At each point, the perceived visual information of the reference and deformed block volumes is obtained by the mutual information  $I(\mathbf{x}|\mathbf{p})$  and  $I(\mathbf{y}|\mathbf{q})$ , respectively. We aim to approximate the perceptual information content from both blocks. To be specific, we subtract the common information

shared between  $\mathbf{p}$  and  $\mathbf{q}$  ( $I(\mathbf{p}|\mathbf{q})$ ) from  $I(\mathbf{x}|\mathbf{q})$  and  $I(\mathbf{y}|\mathbf{q})$ . So, we define a weight based on the mutual information as:

$$w = I(\mathbf{x}|\mathbf{p}) + I(\mathbf{y}|\mathbf{q}) - I(\mathbf{p}|\mathbf{q}), \quad (5)$$

where  $\mathbf{x}$ ,  $\mathbf{y}$ ,  $\mathbf{p}$ , and  $\mathbf{q}$  are all Gaussian for a given  $\alpha$ . The mutual information approximation can be computed using the determinants of the covariances. Details for deriving the mutual information weight function in Equation (5) will be given in the supplementary material.

Let  $x_i$  and  $y_i$  be the  $i$ -th points of two frames in consecutive segments  $\mathbf{X}$  and  $\mathbf{Y}$ , respectively. The Mean Square Error (MSE) between two frames is given by  $\text{MSE} = \frac{1}{P} \sum_{i=1}^P (x_i - y_i)^2$ , where  $P$  is the number of pixels in the frame. For each of the three color channels, we compute a set of weights by moving a sliding voxel across two consecutive segments (see Figure 1), where the spatial size of voxel is  $H \times W$  pixels. This process results in a feature map for each two consecutive segments of the video. We define a weighted MSE for the corresponding location of the central point in the spatial neighborhood using  $w$  defined in Equation (5). Suppose  $x_{j,i}$  and  $y_{j,i}$  are the  $i$ th points at the  $j$ th frame and  $w_{j,i}$  is the weight computed at the corresponding location. The proposed Statistical Spatiotemporal Distillation (SSD) is defined as:

$$\text{SSD}(\mathbf{X}, \mathbf{Y}) = \prod_{j=1}^{\ell} \left( \frac{\sum_{i=1}^P w_{j,i} (x_{j,i} - y_{j,i})^2}{\sum_{i=1}^P w_{j,i}} \right), \quad (6)$$

where  $\ell$  is the length of each segment of the video.

Repeating this process for all two consecutive segments, for a given video of length  $L$ , we can obtain  $\lfloor L/\ell \rfloor - 1$  feature maps per color channel, which encode the appearance and dynamic variations within video segments. This distilled information

can not be used as the input of pretrained CNN models due to multiple channels. To tackle this issue, we use a weighted aggregation approach (*i.e.* weighted sum, for simplicity) to generate a single RGB image from the obtained  $\lfloor L/\ell \rfloor - 1$  channel feature maps. The weights are calculated as:  $\beta_i = \frac{\exp(e_i)}{\sum_i \exp(e_i)}$ , where  $e_i$  is the  $i$ th point of the image map. We compute the weighted sum of points for each channel separately to generate the RGB distilled representation of the video:  $S_j = \sum_{i=1}^{\lfloor L/\ell \rfloor - 1} \beta_i e_{ij}$ , where  $S_j$  denotes the  $j$ th point from one channel of the obtained representation.

Furthermore, we can simplify the mutual information weight function in Equation (5) as:

$$w = \frac{1}{2} \log \left[ \frac{|\mathbf{C}_{(p,q)}|}{\sigma_n^{4K}} \right], \quad (7)$$

where  $K$  is the total number of points in a block volume and

$$|\mathbf{C}_{(p,q)}| = |((\sigma_v^2 + \sigma_n^2)\alpha^2 + \sigma_n^2 g^2 \alpha^2) \mathbf{C}_u + \sigma_n^2 (\sigma_v^2 + \sigma_n^2) \mathbf{I}| \quad (8)$$

Applying an eigenvalue decomposition to the covariance matrix  $\mathbf{C}_u = \mathbf{O} \mathbf{\Lambda} \mathbf{O}^T$ , where  $\mathbf{O}$  is an orthogonal matrix and  $\mathbf{\Lambda}$  is a diagonal matrix with eigenvalues  $\lambda_k$  for  $k = 1, \dots, K$  along its diagonal entries, we can compute  $|\mathbf{C}_{(p,q)}|$  as:

$$|\mathbf{C}_{(p,q)}| = |\mathbf{O} \{(\sigma_v^2 + (1 + g^2)\sigma_n^2)\alpha^2 \mathbf{\Lambda} + \sigma_n^2 (\sigma_v^2 + \sigma_n^2) \mathbf{I}\} \mathbf{O}^T|. \quad (9)$$

Due to the orthogonal property of  $\mathbf{O}$ ,  $|\mathbf{C}_{(p,q)}|$  has a closed-form:

$$|\mathbf{C}_{(p,q)}| = \prod_{k=1}^K \{(\sigma_v^2 + (1 + g^2)\sigma_n^2)\alpha^2 \lambda_k + \sigma_n^2 (\sigma_v^2 + \sigma_n^2)\}. \quad (10)$$

Hence, again, the mutual information weight function in Equation (5) can be expressed as:

$$w = \frac{1}{2} \sum_{k=1}^K \log \left( 1 + \frac{\sigma_v^2}{\sigma_n^2} + \left( \frac{\sigma_v^2}{\sigma_n^2} + \frac{1 + g^2}{\sigma_n^2} \right) \alpha^2 \lambda_k \right). \quad (11)$$

The obtained weight function shows an interesting connection with the local deformation within frames of video. According to the deformation model in Equation (2), the variations from  $\mathbf{x}$  to  $\mathbf{y}$  are characterized by the gain factor  $g$  and the random deformation  $\sigma_v^2$ . As  $g$  is a scale factor along the temporal frame evolution, it does not cause any changes in the image spatial structure. Thus, the changes in image spatial structure are captured by  $\sigma_v^2$ . Our weight function increases monotonically with  $\sigma_v^2$ . This demonstrates that more weights are cast to the areas that have larger variations.

In the derived weight function of Equation (11), we still need to derive approximation for a set of parameters:  $\mathbf{C}_u$ ,  $\alpha^2$ ,  $g$ , and  $\sigma_v^2$ .  $\mathbf{C}_u$  is computed as:

$$\hat{\mathbf{C}}_u = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T, \quad (12)$$

where  $N$  is the number of block volumes and  $\mathbf{x}_i$  is the  $i$ th neighboring vector. The multiplier  $\alpha$  can be approximated using a maximum likelihood estimator:

$$\hat{\alpha}^2 = \frac{1}{K} \mathbf{x}^T \mathbf{C}_u^{-1} \mathbf{x}. \quad (13)$$

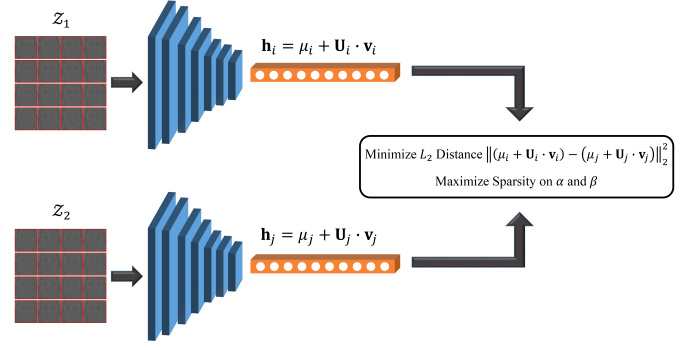


Fig. 2: An illustration of the Siamese network for learning the similarity between samples in an unsupervised manner.

The deformation parameters  $g$  and  $\sigma_v^2$  can be computed via optimizing the following least square regression problem:

$$\hat{g} = \arg \min_g \|\mathbf{y} - g\mathbf{x}\|_2^2, \quad (14)$$

resulting  $\hat{g} = \frac{\mathbf{x}^T \mathbf{y}}{\mathbf{x}^T \mathbf{x}}$ . Substituting this into Equation (2), we can compute  $\sigma_v^2$  using  $\mathbf{v}^T \mathbf{v} / K$ , which leads to:

$$\hat{\sigma}_v^2 = \frac{1}{K} (\mathbf{y}^T \mathbf{y} - \hat{g} \mathbf{x}^T \mathbf{y}). \quad (15)$$

## 2.2. Self-supervised Representation Learning

self-supervised Learning is formulated as firstly learning a pretext task in an *unsupervised* manner, *i.e.* the supervisory labels for the pretext task are automatically generated based on the data itself rather than human annotations, and after the unsupervised training finished, it transfers the learned representations into a downstream supervised task to evaluate the quality of the learned representations. In our problem, we define the pretext task as learning a general similarity function for sample pairs. In the downstream task, we define the downstream task as supervised classification of different pain intensity levels. To encode such function, we propose a method based on Siamese neural networks. We denote a batch of unlabeled data as  $\mathcal{Z} = \{\mathbf{Z}_i\}_{i=1}^N$ , where  $\mathbf{Z}_i$  is a sample instance. Feeding these data into a CNN model  $f(\cdot)$ , we obtain their deep representations from the last fully connected layer. Then, we define a representation matrix  $\mathbf{H} = [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_N]$ , where  $\mathbf{h}_i = f(\mathbf{Z}_i)$  is the extracted deep representation of the  $i$ th sample in the batch.

We propose a joint representation for samples in each batch. The joint representation benefits from constructing a linear model to approximate the structure of each batch in a high-dimensional feature space. We can model a batch as an affine hull of its samples. The affine hull model is used to account for unseen samples in terms of affine combinations of existing samples. Hence, we model a batch as:

$$AH = \left\{ \sum_{i=1}^N \alpha_i \cdot \mathbf{h}_i \mid \sum_{i=1}^N \alpha_i = 1 \right\}, \quad (16)$$

It is possible to represent the affine hull of Equation (16) by another parametric form using batch mean  $\mu$  as a reference point to represent every data:

$$AH = \{\mu + \mathbf{U} \cdot \mathbf{v} \mid \mathbf{v} \in \mathbb{R}^d\}, \quad (17)$$

where the  $d$  columns of  $\mathbf{U}$  are the orthogonal bases obtained from singular value decomposition (SVD) of the centered deep representation matrix. We present a batch of samples as a triplet  $(\boldsymbol{\mu}, \mathbf{U}, \mathbf{H})$  by including both the structure information and deep representations. The information of deep representation is used to reduce the ambiguity of the affine hull space.

Similar samples can be noisy or vulnerable to outliers. Hence, the direct search between nearest neighbors degrades the performance because it is possible to find a pair of similar samples with very small distance that represent different pain intensity levels. To overcome this limitation, we measure the dissimilarity between samples so that the Euclidean distance between them is small and samples of each class could be represented by a combination of its neighborhood points in the deep feature space.

Considering the above point, we propose a convex optimization formulation to define the dissimilarity metric. Given the deep representation matrix of each batch  $\mathbf{H}$ , we evenly split it into two matrices  $\mathbf{H}_i$  and  $\mathbf{H}_j$ . Therefore, their corresponding affine hull representations are  $(\mu_i, \mathbf{U}_i)$  and  $(\mu_j, \mathbf{U}_j)$ . To define the objective function for the convex optimization, we define a series of functions as follows:

$$F_{v_i, v_j} = |(\mu_i + \mathbf{U}_i \cdot \mathbf{v}_i) - (\mu_j + \mathbf{U}_j \cdot \mathbf{v}_j)|_2^2, \quad (18)$$

$$G_{v_i, \alpha} = |(\mu_i + \mathbf{U}_i \cdot \mathbf{v}_i) - \mathbf{H}_i \cdot \alpha|_2^2, \quad (19)$$

$$Q_{v_j, \beta} = |(\mu_j + \mathbf{U}_j \cdot \mathbf{v}_j) - \mathbf{H}_j \cdot \beta|_2^2, \quad (20)$$

where the optimal model coefficients  $\{\mathbf{v}_i^*, \mathbf{v}_j^*\}$  and sample coefficients  $\{\alpha^*, \beta^*\}$  of our metric are achieved by optimizing the following unconstrained objective function:

$$\min_{\mathbf{v}_i, \mathbf{v}_j, \alpha, \beta} F_{v_i, v_j} + \lambda_1 (G_{v_i, \alpha} + Q_{v_j, \beta}) + \lambda_2 \|\alpha\|_1 + \lambda_3 \|\beta\|_1, \quad (21)$$

where the first term is to keep the distance between similar samples  $\mathbf{h}_i = \mu_i + \mathbf{U}_i \cdot \mathbf{v}_i$  and  $\mathbf{h}_j = \mu_j + \mathbf{U}_j \cdot \mathbf{v}_j$  small. The second term is to preserve the individual fidelity between such samples and their neighbors. The last two terms enforces the coefficients to be sparse.  $\lambda_1$ ,  $\lambda_2$ , and  $\lambda_3$  are trade-off values to control the relative importance of different terms.

### 3. Experiments

We evaluated the performance of our proposed method by conducting experiments on two publicly available datasets, namely UNBC-McMaster Shoulder Pain Expression Archive (Lucey et al. (2011)) and the BioVid Heat Pain - Part A (Walter et al. (2013)). Table 1 summarizes the content of these datasets.

#### 3.1. Experimental Setup

In our experiments, we use four different network architectures, *i.e.* VGG16 (Simonyan and Zisserman (2014)), Inception-V1 (Ioffe and Szegedy (2015)), ResNet-18 (He et al. (2016)), and ResNet-50 (He et al. (2016)). All of these models are trained using stochastic gradient descent with the momentum of 0.9 and an annealed learning rate, starting from

Table 1: Characteristics of the pain recognition datasets.

Dataset	Subjects	Videos	Stimuli	Classes
UNBC-McMaster	25	200	Shoulder Pain	16
BioVid (Part A)	87	8,700	Induced Heat	5

$10^{-3}$  and multiplied by a factor of 0.2 per epoch. During the training, we randomly perform size jittering, cropping, flipping, and rescaling on the input samples. We also applied our SSD on optical flow data. For the computation of the optical flow, we use TLV1 optical flow algorithm (Zach et al. (2007)), which is implemented in OpenCV with CUDA.

The value of  $\lambda_1$  is fixed as 0.01 for all the experiments conducted in this paper. For  $\lambda_2$  and  $\lambda_3$ , we design an automatic mechanism to control the relative sparsity of  $\alpha$  and  $\beta$ . Note that if  $\lambda_2 \geq \max(|2\lambda_1 \cdot (\mathbf{H}^T \boldsymbol{\mu})|)$ , the zero vector is optimal for  $\alpha$  at zero. The same statement applies for  $\lambda_3$  and  $\beta$ . By performing a grid search and following the guidelines in (Bengio (2012)), we adaptively set  $\lambda_2 = 0.1\lambda_2^*$  and  $\lambda_3 = 0.1\lambda_3^*$  for all experiments.

#### 3.2. Analysis of SSD Representations

In this section, we evaluate the performance of our method by analyzing different parameters in SSD representations. As described in Section 2.1, SSD divides the video into smaller fix-sized non-overlapping segments. The spatiotemporal patch-wise statistical analysis between two consecutive segments allows SSD to capture and encode the appearance and dynamic variations into one single RGB image map. Hence, the performance of SSD depends on two parameters, *i.e.* the length of segments and the spatial neighborhood around each point. We analyze the performance of our method by varying the length of video segments from 5 to 60 frames. Figure 3 shows the results of experiments using the UNBC-McMaster dataset. As depicted, the Area Under the Curve (AUC) accuracy increases as the number of frames per segment increases. However, after a certain number of frames (20 frames per segment), the performance drops significantly. The deterioration in performance is likely due to capturing a wide range of temporal information from the face, which covers different pain intensity levels throughout the segment. Hence, SSD is not able to effectively encode the appearance and dynamic of the sequence.

Another important parameter, which influences the quality of obtained representations using SSD, is the spatial size of the block volumes. We change the spatial size of neighborhood around each point of the video segments to investigate its effect. Figure 3 demonstrates the results of this experiment. Using the small spatial windows, the receptive field of SSD is also small. Hence, the statistical operation is performed in a tiny region of the video, which ignores most important relationships between neighboring points. However, choosing a large size for spatial patches leads to low accuracy. We argue that this drop in the performance is due capturing too much information and increasing the complexity of the obtained image map. Based on the results of this experiment, we will use video segments of 20 frames with the spatial size of  $5 \times 5$  in the rest of our experiment, unless otherwise it is mentioned.

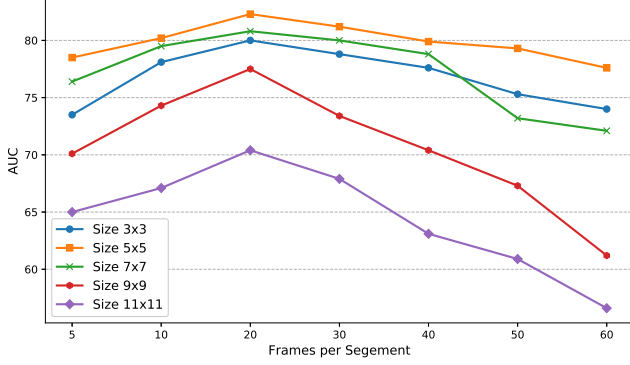


Fig. 3: Accuracy (%) of the proposed method by varying the segment length size and the spatial size of neighborhood for SSD representation computation using ResNet-50.

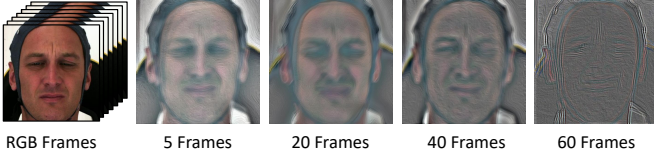


Fig. 4: Visualization of SSD representations versus different segment lengths. Small segments do not cover wide range of spatiotemporal variations, while too long segments have more noise.

In order to gain a better insight into the effect of segment length on SSD, we computed SSD representations using different segment lengths in Figure 4. We varied the number of frames per segment from 5 to 60 frames. As can be seen, the smaller segment, lower spatiotemporal variations. In contrast, SSD representations obtained from longer segments (e.g. 60 frames) show noisy behavior, which interferes the distilled spatiotemporal information. Based on Equations (1) and (2), enlarging the temporal size of segments increases the spatiotemporal information to noise ratio, while decreasing the temporal size of segments limits SSD computation to a small portion of the video. Hence, it is important to select sufficiently enough number of frames per segments for SSD computation. This qualitative analysis of SSD representations further validates the quantitative results in Figure 3.

We compared our proposed SSD representation against two frame-based baseline models. The first model operates on RGB images and the second model works on the optical flow. We also applied SSD on the optical flow sequence and the RGB video. The results are summarized in Table 2. The proposed SSD has 2.3% and 1.2% improvement over optical flow on the UNBC-McMaster and the BioVid datasets, respectively, indicating the superiority of SSD in capturing the appearance and dynamic of the video.

### 3.3. 2D Models vs. 3D Models

The proposed SSD does not require any parameter learning process. It captures appearance and dynamic information of the video and encodes it into one RGB image. This enables its output representations to be processed by models devised for images. Consequently, the need for complex 3D models

Table 2: The AUC accuracy (%) of our proposed method versus the accuracy of representations from optical flow and RGB data using ResNet-50.

Data Type	UNBC-McMaster	BioVid
RGB	85.9	69.4
OF	86.2	69.8
SSD (RGB)	88.5	71.0
SSD (OF)	86.6	70.3

Table 3: Average testing runtime (sec.) per frame and accuracy (%) on the UNBC-McMaster.

Architecture	2D		3D	
Method	Time	AUC	Time	AUC
VGG16	0.924	75.3	2.769	77.1
Inception-V1	0.985	78.5	2.854	78.8
ResNet-18	1.085	80.5	2.866	79.3
ResNet-50	1.299	83.2	2.892	81.5

and large amount of annotated training videos can be avoided. Although extracting the distilled representations adds overhead to training of 2D models, this time is negligible when compared with the training time of 3D models. For instance, the training time of ResNet-50 on RGB frames of the UNBC-McMaster is 1,257 minutes, which is raised to 1,326 minutes by applying SSD. However, our method is more efficient in the test time. To demonstrate this efficiency, we randomly selected a set of 50 videos from the test set of the UNBC-McMaster dataset. The videos have lengths in range of 150 – 200 frames. We, then, compare the runtime of our method using 2D models against their 3D counterparts in Table 3, where we replace 2D filters of CNNs by 3D filters. From Table 3, we conclude that encoding information via SSD enables us to estimation pain intensity levels by using less complex models devised for images and still achieve high performance.

### 3.4. Cross Database Analysis

The performance analysis of automatic pain assessment methods in the cross-database setting is undermined. The existing approaches' performance degrades when the training and testing datasets are different (Othman et al. (2019)). This is mainly due to lack of generalization in the learned representations. In this section, we conduct cross-dataset experiments to show the generalization of the learned representations in our self-supervised framework. In these experiments, the model learns the pretext task using unlabeled data (SSD representations of videos) of either the UNBC-McMaster, BioVid, or their combination and is fine-tuned and tested using labeled SSD representations of the other dataset. We combined different portion of the target data with the training data for learning the pretext task to analyze how the performance is influenced using combined datasets. Table 4 summarizes the results. As can be seen, the model achieves reasonably good performance in self-supervised fashion and benefits from a larger dataset in the pretext task learning. However, including a portion of target

Table 4: The accuracy (%) of the proposed self-supervised learning method in cross-dataset setting. ResNet-50 is used as the backbone of our model.

Pretext Data	Downstream Data	AUC
BioVid	UNBC-McMaster	69.2
UNB-McMaster	BioVid	65.5
BioVid + UNBC-McMaster (10%)	UNBC-McMaster	71.3
BioVid + UNBC-McMaster (30%)	UNBC-McMaster	73.5
BioVid + UNBC-McMaster (50%)	UNBC-McMaster	74.2
UNBC-McMaster + BioVid (10%)	BioVid	64.4
UNBC-McMaster + BioVid (30%)	BioVid	69.0
UNBC-McMaster + BioVid (50%)	BioVid	71.8

Table 5: Comparison of supervised pain intensity estimation using SSD representations in the cross-dataset experimental setting using ResNet-50.

Pretext Data	Downstream Data	AUC
BioVid	UNBC-McMaster	80.1
UNB-McMaster	BioVid	75.5
BioVid + UNBC-McMaster (10%)	UNBC-McMaster	81.2
BioVid + UNBC-McMaster (30%)	UNBC-McMaster	82.3
BioVid + UNBC-McMaster (50%)	UNBC-McMaster	83.5
UNBC-McMaster + BioVid (10%)	BioVid	76.3
UNBC-McMaster + BioVid (30%)	BioVid	78.2
UNBC-McMaster + BioVid (50%)	BioVid	79.0

data to the training data of the pretext task can slightly improve the performance.

To gain a better insight in effectiveness of our proposed self-supervised method, we conduct complementary cross-dataset experiments in a supervised manner. In these experiments, we train the CNN model using SSD representations of one dataset (with labels) and, then, use the SSD representations of other dataset for testing. Table 5 lists the results of this experiment. Although the accuracy drops in the cross-dataset setting, the SSD still shows a good performance, demonstrating its high capability to represent video sequences discriminatively. From the results of Tables 4 and 5, we notice that there is still a gap in performance between supervised and self-supervised techniques. However, the cross-dataset experiments show promising results by using SSD representations.

### 3.5. Comparison against the State-of-the-Art

We compare the performance of our proposed method with the-state-of-the-art methods for pain intensity estimation on the UNBC-McMaster (Lucey et al. (2011)) and the BioVid (Walter et al. (2013)) datasets. In these experiments, VGGFace2 dataset (Cao et al. (2018)) is used as unlabeled data for learning the pretext task. We conduct experiments using the original implementation of supervised and self-supervised approaches on both RGB and SSD representations. To make a direct and fair comparison, we report the Mean Square Error (MSE) and the Pearson Correlation Coefficient (PCC) for the UNBC-McMaster dataset. Table 6 summarizes the comparative results on the UNBC-McMaster dataset. We observe that the proposed method improves the performance of self-supervised pain intensity estimation using either RGB

Table 6: Comparative analysis in terms of the Mean Square Error (MSE) and the Pearson Correlation Coefficient (PCC) on the UNBC-McMaster dataset using RGB and SSD representations. ResNet-50 is used as the backbone of our method.

	Method	RGB		SSD	
		MSE	PCC	MSE	PCC
Supervised	RVR (Kaltwang et al. (2012))	1.39	0.59	1.15	0.63
	HoT (Florea et al. (2014))	1.21	0.53	0.98	0.66
	OSVR (Zhao et al. (2016))	N/A	0.60	1.06	0.63
	RCNN (Zhou et al. (2016))	1.54	0.64	1.27	0.73
	LSTM (Rodriguez et al. (2017))	0.74	0.78	0.69	0.80
	SCN (Tavakolian and Hadid (2019))	0.32	0.92	N/A	N/A
Self-supervised	Geometry (Gan et al. (2018))	1.81	0.51	1.65	0.57
	OPN (Lee et al. (2017))	1.76	0.57	1.60	0.61
	Jigsaw (Noroozi and Favaro (2016))	1.73	0.68	1.58	0.70
	Wang et al. (2019)	1.59	0.65	1.33	0.70
	Caron et al. (2018)	1.47	0.71	1.12	0.73
	Our Method	1.03	0.74	0.92	0.78

Table 7: Comparative analysis in terms of the AUC accuracy (%) on the BioVid dataset. ResNet-50 is used as the backbone of our method.

	Method	RGB	SSD
		AUC	AUC
Supervised	Head-movement (Walter et al. (2013))	67.00	70.36
	Time-windows (Walter et al. (2013))	71.00	72.15
	LBP (Yang et al. (2016))	63.72	65.20
	BSIF (Yang et al. (2016))	65.17	67.83
	FAD set (Werner et al. (2017))	72.40	74.59
	SCN (Tavakolian and Hadid (2019))	86.02	N/A
Self-supervised	Geometry (Gan et al. (2018))	61.29	62.50
	OPN (Lee et al. (2017))	63.07	64.83
	Jigsaw (Noroozi and Favaro (2016))	63.44	65.19
	Wang et al. (2019)	65.23	68.41
	Caron et al. (2018)	66.01	69.19
	Our Method	69.35	71.02

or SSD data following leave-one-subject-out cross-validation. Our method achieve 1.03 and 0.92 MSE using RGB and SSD data, respectively. It also improves PCC among self-supervised methods. We assert that this improvement is due to learning generalized representations in the pretext task thanks to the proposed similarity function, which push dissimilar samples away from each other. However, there is still a margin between the performance of the self-supervised and the supervised approaches in Table 6.

Table 7 draws comparison against the state-of-the-art methods on the BioVid dataset. We use part A of this dataset, which contains only unoccluded facial videos. As can be seen, our method achieves 69.35% and 71.02% AUC accuracy using RGB and SSD data, respectively, and improves the highest performance of self-supervised techniques by 3.34% and 1.83%, showing results comparable to some of the supervised techniques.

## 4. Conclusion

In this paper, we proposed a novel self-supervised representation learning framework for automatic pain intensity

estimation. By learning a similarity function on a large portion of unlabeled samples using a Siamese network as the pretext task, we achieved generalized representations of data in an unsupervised manner. The learned representations were further transferred to the downstream supervised task, where the CNN model is fine-tuned using a small subset of labeled samples. In order to reduce the complexity of the learning process, we presented a statistical spatiotemporal distillation technique to capture and encode the appearance and dynamic of the video into one single RGB image. Hence, we could use 2D models to process video data. We conducted extensive experiments on two publicly available datasets to demonstrate the effectiveness of our proposed method. The experimental results of the cross-dataset pain intensity estimation validated our initial hypothesis of the generalization of the learned representations.

## Acknowledgments

We would like to acknowledge the financial support of the Academy of Finland (No. 331883), Infotech Oulu, Tauno Tönning, Nokia, and KAUTE foundations.

## References

- Bengio, Y., 2012. Practical recommendations for gradient-based training of deep architectures, in: *Neural networks: Tricks of the trade*. Springer, pp. 437–478.
- Bromley, J., Guyon, I., LeCun, Y., Säckinger, E., Shah, R., 1994. Signature verification using a "siamese" time delay neural network, in: *NIPS*, pp. 737–744.
- Cao, Q., Shen, L., Xie, W., Parkhi, O.M., Zisserman, A., 2018. VGGFace2: A dataset for recognising faces across pose and age, in: *FG*.
- Caron, M., Bojanowski, P., Joulin, A., Douze, M., 2018. Deep clustering for unsupervised learning of visual features, in: *ECCV*, pp. 132–149.
- Dailey, M.N., Joyce, C., Lyons, M.J., Kamachi, M., Ishi, H., Gyoba, J., Cottrell, G.W., 2010. Evidence and a computational explanation of cultural differences in facial expression recognition. *Emotion* 10, 874–893.
- Ekman, P., Friesen, W.V., 1978. *Facial Action Coding System: A Technique for the Measurement of Facial Movements*. Consulting Psychologist Press.
- Florea, C., Florea, L., Vertan, C., 2014. Learning pain from emotion: Transferred HoT data representation for pain intensity estimation, in: *ECCV Workshops*, pp. 778–790.
- Gan, C., Gong, B., Liu, K., Su, H., Guibas, L.J., 2018. Geometry guided convolutional neural networks for self-supervised video representation learning, in: *CVPR*, pp. 5589–5597.
- He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition, in: *CVPR*, pp. 770–778.
- Hossein Khatonabadi, S., Vasconcelos, N., Bajic, I.V., Shan, Y., 2015. How many bits does it take for a stimulus to be salient?, in: *CVPR*, pp. 5501–5510.
- Hu, J., Lu, J., Tan, Y.P., 2014. Discriminative deep metric learning for face verification in the wild, in: *IEEE CVPR*, pp. 1875–1882.
- Ioffe, S., Szegedy, C., 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*.
- Jing, L., Tian, Y., 2020. Self-supervised visual feature learning with deep neural networks: A survey. *IEEE PAMI*.
- Kaltwang, S., Rudovic, O., Pantic, M., 2012. Continuous pain intensity estimation from facial expressions, in: *Advances in Visual Computing*, pp. 368–377.
- Lee, H.Y., Huang, J.B., Singh, M., Yang, M.H., 2017. Unsupervised representation learning by sorting sequences, in: *ICCV*, pp. 667–676.
- Littlewort, G.C., Bartlett, M.S., Lee, K., 2009. Automatic coding of facial expressions displayed during posed and genuine pain. *Image and Vision Computing* 27, 1797–1803.
- Lucey, P., Cohn, J.F., Prkachin, K.M., Solomon, P.E., Chew, S., Matthews, I., 2012. Painful monitoring: Automatic pain monitoring using the UNBC-McMaster shoulder pain expression archive database. *Image and Vision Computing* 30, 197–205.
- Lucey, P., Cohn, J.F., Prkachin, K.M., Solomon, P.E., Matthews, I., 2011. Painful data: The UNBC-McMaster shoulder pain expression archive database, in: *FG*.
- Noroozi, M., Favaro, P., 2016. Unsupervised learning of visual representations by solving jigsaw puzzles, in: *ECCV*, pp. 69–84.
- Othman, E., Werner, P., Saxen, F., Al-Hamadi, A., Walter, S., 2019. Cross-database evaluation of pain recognition from facial video, in: *ISPA, IEEE*, pp. 181–186.
- Rodriguez, P., Cucurull, G., Gonzalez, J., Gonfaus, J.M., Nasrollahi, K., Moeslund, T.B., Roca, F.X., 2017. Deep pain: Exploiting long short-term memory networks for facial expression classification. *IEEE Trans. Cybernetics*, 1–11.
- Simonyan, K., Zisserman, A., 2014. Very deep convolutional networks for large-scale image recognition. *arXiv:1409.1556*.
- Tavakolian, M., Hadid, A., 2019. A spatiotemporal convolutional neural network for automatic pain intensity estimation from facial dynamics. *IJCV* 127, 1413–1425.
- Walter, S., Gruss, S., Ehleiter, H., Tan, J., Traue, H.C., Crawcour, S., Werner, P., Al-Hamadi, A., Andrade, A.O., 2013. The biovid heat pain database data for the advancement and systematic validation of an automated pain recognition system, in: *Int. Conf. Cybernetics*.
- Wang, J., Jiao, J., Bao, L., He, S., Liu, Y., Liu, W., 2019. Self-supervised spatio-temporal representation learning for videos by predicting motion and appearance statistics, in: *CVPR*, pp. 4006–4015.
- Wang, L., Xiong, Y., Wang, Z., Qiao, Y., Lin, D., Tang, X., Van Gool, L., 2018a. Temporal segment networks for action recognition in videos. *IEEE Trans. PAMI*.
- Wang, X., Wang, X., Ni, Y., 2018b. Unsupervised domain adaptation for facial expression recognition using generative adversarial networks. *Computational intelligence and neuroscience* 2018, 1–10.
- Werner, P., Al-Hamadi, A., Limbrecht-Ecklundt, K., Walter, S., Gruss, S., Traue, H.C., 2017. Automatic pain assessment with facial activity descriptors. *IEEE Trans. Affective Computing* 8, 286–299.
- Yan, H., 2016. Transfer subspace learning for cross-dataset facial expression recognition. *Neurocomputing* 208, 165–173.
- Yang, R., Tong, S., Boddallo, M., Boutellaa, E., Peng, J., Feng, X., Hadid, A., 2016. On pain assessment from facial videos using spatio-temporal local descriptors, in: *IPTA*.
- Zach, C., Pock, T., Bischof, H., 2007. A duality based approach for realtime tv-l 1 optical flow, in: *Joint pattern recognition symposium*, pp. 214–223.
- Zhao, R., Gan, Q., Wang, S., Ji, Q., 2016. Facial expression intensity estimation using ordinal information, in: *CVPR*, pp. 3466–3474.
- Zhou, J., Hong, X., Su, F., Zhao, G., 2016. Recurrent convolutional neural network regression for continuous pain intensity estimation in video, in: *CVPR Workshops*, pp. 84–92.