

# Sparsity in Dynamics of Spontaneous Subtle Emotions: Analysis & Application

Anh Cat Le Ngo, *Member, IEEE*, John See, *Member, IEEE*, Raphael C.-W. Phan, *Member, IEEE*

**Abstract**—Subtle emotions are present in diverse real-life situations: in hostile environments, enemies and/or spies maliciously conceal their emotions as part of their deception; in life-threatening situations, victims under duress have no choice but to withhold their real feelings; in the medical scene, patients with psychological conditions such as depression could either be intentionally or subconsciously suppressing their anguish from loved ones. Under such circumstances, it is often crucial that these subtle emotions are recognized before it is too late. These spontaneous subtle emotions are typically expressed through micro-expressions, which are tiny, sudden and short-lived dynamics of facial muscles; thus, such micro-expressions pose a great challenge for visual recognition. The abrupt but significant dynamics for the recognition task are temporally sparse while the rest, i.e. irrelevant dynamics, are temporally redundant. In this work, we analyze and enforce sparsity constraints to learn significant temporal and spectral structures while eliminating irrelevant facial dynamics of micro-expressions, which would ease the challenge in the visual recognition of spontaneous subtle emotions. The hypothesis is confirmed through experimental results of automatic spontaneous subtle emotion recognition with several sparsity levels on CASME II and SMIC, the two well-established and publicly available spontaneous subtle emotion databases. The overall performances of the automatic subtle emotion recognition are boosted when only significant dynamics of the original sequences are preserved.

**Index Terms**—Spontaneous subtle emotions, emotion suppression, data sparsity, dynamic mode decomposition, micro-expression recognition.



## 1 INTRODUCTION

In our current era of social networks facilitated primarily by the widespread accessibility of internet-ready on-person mobile devices, smart human-centric systems are increasingly expected to perceive and understand humans rather than vice versa. Instead of simply executing users' commands, computers need to understand the multi-modality of human-like communications. This includes recognition of facial expressions, the non-verbal form of human communications. The shift in paradigm toward human-friendly computing has initiated the field of Affective Computing [1]. This section briefly surveys recent methods and advances in an emerging subfield of Affective Computing: automatic recognition of subtle emotions through facial micro-expressions. Though automatic recognition of spontaneous subtle emotions is a new challenging task, recognition systems for normal expressions have been a research subject for nearly two decades. Shan et al. [2] summarized achievements of this research field in a popular framework of automatic facial expression and emotion recognition systems. Furthermore, the authors focused on analyzing two core components: facial representations e.g. LBPTOP, and classifiers e.g. SVM and AdaBoost. For normal expressions, this framework successfully recognizes with above 90% accuracy [2]. However, the same framework falls short of this impressive recognition rate when dealing with micro-expressions [3], [4].

As micro-expressions of subtle emotions are much more elusive than normal expressions due to their small intensities, short-liveness (between  $\frac{1}{25}$ s and  $\frac{1}{15}$ s [5]) and unpredictability, their image sequences need to be pre-processed to reduce these unfavorable characteristics. In this paper, we firstly propose removal of redundant neutral faces from micro-expression sequences and keeping only sparse and significant frames, which is illustrated in Figure 1. As the sparse frames significantly contribute to reconstruction of original dynamics; meanwhile, redundant frames could be omitted without much cost i.e. errors between reconstructed and original sequences. The proposed pre-processing technique aims to remove as many neutral and redundant frames as possible with minimum cost. It not only produces more visually distinguishable but also allows extraction of more discriminant features. As a result, it improves the accuracy rate of automatic subtle emotion recognition. Secondly, we carry out temporal and spectral analysis of subtle emotion sequences so as to clarify rationales behind our approaches as well as select suitable experimental parameters. Finally, we compare performances of our proposed solution with those of the state-of-the-art methods on recognition of spontaneous micro-expressions.

Section 1.1 describes two publicly available spontaneous subtle emotions databases: CASME II [3] and SMIC [6], which are utilized as input data for our experiments. As system performance greatly depends on characteristics of these databases, understanding their pros and cons as well as samples provides knowledge and clues for designing optimal automatic recognition systems. In Section 1.2, related works in recent literatures are reviewed and categorized with respect to their main contributions in preprocessing, feature extraction or classification stages. Section 2 elabo-

Anh Cat Le Ngo and Raphael C.-W. Phan are with the Faculty of Engineering, Multimedia University, Malaysia.  
John See is with the Faculty of Computing & Informatics, Multimedia University, Malaysia.

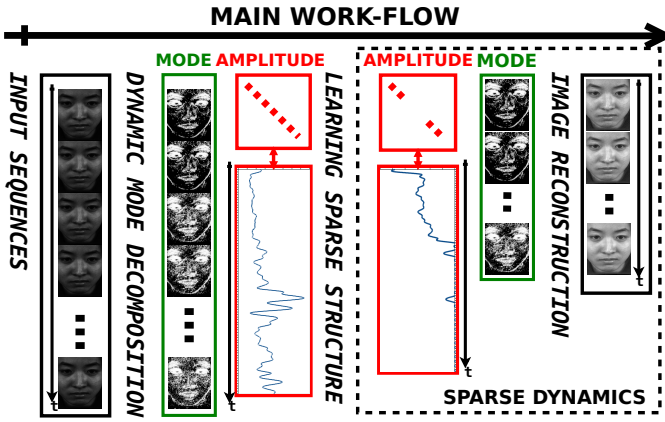


Fig. 1: Visualization of steps in Sparse Promoting Dynamic Mode Decomposition (DMDSP) for removals of redundant neutral faces: (1) Input sequences are analyzed by Dynamic Mode Decomposition (DMD), (2)  $L_1$  sparsity is applied to achieve the least number of modes and give the minimum loss during reconstruction, (3) Output sequences are reconstructed from sparse modes.

rates details of dynamic preprocessing techniques: Temporal Interpolation Method (TIM) [6], Dynamic Mode Decomposition (DMD) [7] and Sparsity-promoting Dynamic Mode Decomposition (DMDSP) [8]. Meanwhile, Section 3 utilizes DMD magnitudes to analyze responses of temporally dense (TIM) and sparse (DMDSP) sampling approaches in the temporal and spectral domains. Section 4 describes with what parameters and how the proposed method are evaluated as well as compared to both baselines [3] [9] and other recent state-of-the-art methods [10], [11], [12]. Finally, conclusions are drawn in Section 5.

### 1.1 Spontaneous Subtle Emotion Databases

Various facial expression databases had been proposed in the literature. However, little attention has been paid to *spontaneous* micro-expression databases, partly due to difficulties faced in proper elicitation of samples and labeling by experts. Hence, the lack of well-established databases for spontaneous micro-expression researches poses a challenge towards the design of automatic detection and recognition systems. There is a need to emphasize on the actual meaning of the term “spontaneous micro-expressions”, in contrast to what had been previously regarded generally as “micro-expressions”. Ekman’s [5] and Yan et al. [3] suggest that micro-expressions should be considered involuntary and difficult to disguise. Previous micro-expression databases such as the USF-HD [13] and Polikovskiy’s [14] databases contained micro-expressions that were actually posed or acted out instead of naturally spontaneous ones. Moreover, the occurrence duration of their micro-expressions were longer ( $\frac{2}{3}$ s) than Paul Ekman’s definition ( $\frac{1}{3}$ s) [5]. Another database is the YorkDDT [15], which includes the micro-expression in a spontaneous manner. However, there were other irrelevant facial and head movements therein, thus complicating the recognition process. Besides, YorkDDT contained only 18 micro-expressions which are insufficient for proper experimentation and analysis.

However, there are two recent and comprehensive databases that meet the requirements of spontaneous micro-expressions, namely the SMIC<sup>1</sup> [6] and CASME II<sup>2</sup> [16] (an improved extended version of the original CASME by the same researchers). CASME II has 247 video samples from 26 subjects and SMIC has 164 samples from 16 subjects. They are both publicly available and contain sufficiently large number of video samples which are conducive for a micro-expression recognition research. Both databases were recorded in a constrained laboratory condition annotated by two trained coders and also the participants’ self-reports. Non-emotional facial movements were also eliminated from the final selected sequence. Samples from both SMIC and CASME II were acquired from relatively high frame rates (100 fps and 200 fps respectively) to better locate the occurrence of micro-expressions.

### 1.2 Related Work

Discussion about recent updates in the three main stages of the subtle emotion recognition system: pre-processing, facial representation and classification, is presented in the following subsections 1.2.1, 1.2.2, and 1.2.3 respectively.

#### 1.2.1 Subtle Emotion Preprocessing

The subtleness of spontaneous emotions is challenging to be recognized due to two main problems: small dynamics of facial muscles, and involuntary and unexpected expressions. Therefore, video samples need pre-processing steps to better visualize changes in subtle emotions and subsequently extract more distinctive features. As motions of facial muscles in micro-expressions are too small, Le Ngo et al. [17] showed that motion magnification techniques [18] [19] [20] improve the recognizability of these expressions. These techniques are able to increase the emotional intensity of micro-expressions, making them more visible like normal expressions. Moreover, this magnification effect can be achieved by fast Eulerian Motion Magnification techniques [18] [19] instead of the Lagrangian approach [20] which often requires motion estimations, and other heavily computational processes. However, magnification of micro-expression is outside the scope of this work as this paper mainly focuses on the sparsity of these expressions.

While motion magnification deals with small displacements between frames of micro-expression video samples, temporal interpolation method (TIM) [6] tackles the unexpectedness of micro-expressions. Bursts of spontaneous subtle emotions are difficult to be detected accurately; therefore, video samples are often cut from a long recording of a subject’s expressions. While on-set and off-set points are identified by trained experts in subtle emotions to indicate the starting and ending of a micro-expression sequence, these points are hardly accurate as well. Therefore, these video samples may include frames of almost neutral faces among frames of micro-expressions. Since micro-expressions only last for a very short duration, neutral faces may dominate a large portion of some sequences. For a sample with many redundant neutral faces, TIM is able to interpolate at arbitrary points along a temporal axis according to an

1. <http://www.cse.oulu.fi/SMICDatabase>  
 2. <http://fu.psych.ac.cn/CASME/casme2-en.php>

embedded graph in a manifold, which is in turn learned from video frames. TIM was initially aimed at synthesizing more frames, as video samples recorded at standard 25 fps were too short for subtle emotion recognition. However, the same technique could also be used to interpolate less frames or to remove redundant neutral faces, as video samples were recorded at 100 fps or 200 fps are too long. As TIM assumed that facial expressions change across consecutive frames, and are sampled along a simple graph on a manifold, it is difficult to control how significant or redundant the dynamics are after the interpolation. Therefore, the positive effectiveness of TIM on the performance of the recognition system cannot be guaranteed.

A technique capable of extracting coherent structures and significant dynamics at a single temporal frequency, is Dynamic Mode Decomposition (DMD) [7]. DMD is a popular technique in fluid dynamics imagery, and it was recently applied to foreground motion segmentation in video processing [21]. A more recent variant of it, Sparsity-Promoting Dynamic Mode Decomposition (DMDSP) [8] puts the decomposition under sparse constraints such that the least number of DMD modes are utilized for construction of original sequences. The notion of analyzing a sequence of images into more meaningful temporal structures is potentially useful; an idea which we aim to exploit in this paper.

### 1.2.2 Subtle Emotion Features

Systems for automatic recognition or detection of micro-expressions inherited many components from those for normal expressions (or so-called macro-expressions), including use of features. As Local Binary Pattern with Three Orthogonal Planes (LBPTOP) [22] is a common and effective feature for representing normal facial expressions [2], it has also been utilized in several works relating to micro-expressions [3], [4], [6]. LBPTOP is a spatiotemporal feature, encoding textural features along three orthogonal physical planes  $XY$ ,  $YT$ , and  $XT$  into binary sequences, where  $X, Y$  are two axes of the spatial domain and  $T$  is the temporal axis. The binary sequences are later summarized and concatenated in a histogram, which forms the LBPTOP feature. Local Spatiotemporal Directional Features (LSDF) was recently proposed by Wang et al. [16] for automatic recognition of subtle emotions. Instead of using the center pixel of the neighborhood for thresholding, LSDF encodes each plane along the horizontal and vertical directions. Their experiments demonstrated that LSDF was comparable to LBPTOP in most cases, if not better under certain conditions.

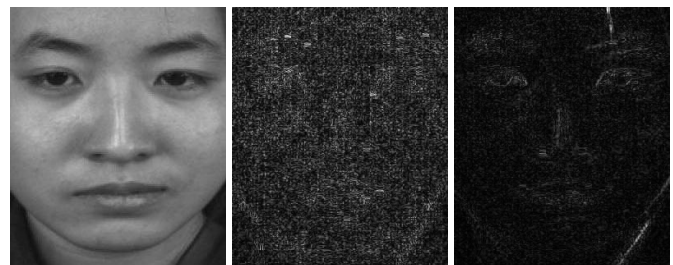
Besides statistical spatiotemporal textural features e.g. LBPTOP, there are other potential feature extraction and representations for micro-expressions based on optical flow, multi-scale wavelet analysis, etc. For instance, Liong et al. [12] utilized optical strain, a derivative of optical flow, for the recognition task; building on similar concepts used for expression spotting by Shreve et al. [13], [23]. Furthermore, Liu et al. [24] encoded statistical information of the main directional optical flows in regions of interests (ROI), which are manually defined with respect to facial landmarks. A recent work by Oh et al. [11] introduced a multi-scale Riesz wavelet representation for micro-expressions that captured the monogenic signal components, i.e. magnitude, phase and orientation. Their method reported an improvement

over the LBPTOP and spatiotemporal local monogenic binary pattern (STLMBP) [25]. Furthermore, Oh et al. [26] proved experimentally that intrinsic 2-D features are better than its 1-D counterpart for encoding facial micro-expressions. Recently, Huang et al. [10] have proposed a new spatio-temporal feature based on integral projection of difference images along horizontal and vertical directions, which achieved a good recognition performance relative to current state-of-the-art methods. Despite different approaches, all features aim to extract error-prone dynamics and tiny motions of micro-expressions, which are discriminative features for recognizing subtle emotions.

### 1.2.3 Subtle Emotion Classifier

Besides features, choices of classifiers for micro-expression detection and recognition are inherited from approaches for normal macro-expressions. Shan et al. [2] highlighted two popular classifiers for emotion recognition systems: Support Vector Machine SVM and AdaBoost, both of which have also been utilized for recognizing subtle emotions [3], [4]. However, there is a distinct difference between the distribution of samples in macro- and micro-expression databases, which directly affects performance and choices of classifiers. Since video samples of spontaneous macro-expressions are widely available with a large number of samples, it is easy to get a balance between classes of macro-expressions. It is much more difficult to acquire a balanced number of spontaneous micro-expression samples due to various reasons, viz. its natural characteristic (small intensities, unexpected dynamics) and the difficulty in eliciting certain emotions. Therefore, imbalance of samples across classes is unavoidable. To tackle this imbalance, Le Ngo et al. [4] introduced an Adaboost-based person-specific classifier, and advocated the use of F1-score, precision and recall metrics in place of the conventional recognition accuracy. Their experimental results showed an improvement by a small margin when compared to standard classifiers like SVM and AdaBoost.

## 2 DYNAMICALLY PREPROCESSING METHODS



(a) Subtle Emotion (b) DMD Mode (c) DMDSPP Mode

Fig. 2: Visual comparison between a noisy and redundant DMD spatial mode in (b) and a clear and significant DMDSPP spatial mode in (c) from a subtle expression in (a)

Most subtle emotions happen very briefly in a short period of time; therefore only high-speed recording is able to capture their full dynamics. Moreover, these expressions usually appear unexpectedly and their (beginning) on-set and (ending) off-set points are difficult to be identified

exactly even by trained experts. As a result, more unnecessary frames, containing no emotional expressions, are accidentally acquired for a micro-expression sample. The redundancy is inevitable, as shown in the Figure 2 visualization of redundant and significant spatial modes of a subtle expression. Due to the inseparability of identities and emotions, unnecessary neutral faces only confuse classifiers of subtle expressions and dampen performances in the latter recognition task. Hence, removal of these undesired frames is crucial.

Lets consider a discrete signal, in Figure 3a, which represents dynamic magnitudes  $f(t)$  of a spontaneous subtle emotion at time  $t$ . This toy example deliberately demonstrates the redundancy assumption of the dynamics  $f(t)$  as only two short parts of the sample signal have significant magnitudes while the majority of these discrete signals have relatively small magnitudes. In other words, significant dynamics are sparse and insignificant ones are redundant for reconstruction of facial dynamics as most signal energy is concentrated into these two local peaks. We hypothesize that micro-expressions would become more descriptive and discriminative if only significant dynamics are kept and redundant ones are eliminated. In the following Subsections 2.1 and 2.2, we discuss two approaches to deal with this redundancy – Uniform Sampling and Sparse Sampling approaches; the latter being our proposed scheme.

## 2.1 Uniform Sampling Approach: Temporal Interpolation Method (TIM)

In this approach, micro-expressions are assumed to happen continuously and sampled along a curve on a low-dimensional manifold. If the curve and its manifold are successfully parameterized, a specific number of frames can be synthesized or interpolated from the original video frames. Figure 3b demonstrates how this so-called *uniform selection* approach fits the original discrete signals in a curve (dashed line) and synthesizes samples at arbitrary but equispaced points along that curve (circular markers). The synthesized discrete signal can either interpolate towards a reduced number of samples (as seen in Figure 3b), or extrapolate to more samples if necessary.

One example of an approach that performs the described uniform selection is the Temporal Interpolation Method (TIM), used by Pfister et al. [6] and first suggested by Zhou et al. [27]. While Zhou et al. extrapolated frames for practical lip-reading, Pfister et al. aimed for recognizing subtle emotions from normal frame-rate recorded sequences. Both applications aimed to synthesize proper frame-lengths for stable spatio-temporal feature extraction. TIM was utilized by Pfister et al. in the opposite way (interpolation) instead of extrapolation. As spontaneous subtle emotions are recorded with high frame-rate, TIM interpolates fewer frames than original frame-lengths to remove redundancy in the dynamics of micro-expressions. Interpolated frames are uniformly sampled at equispaced positions of a graph embedded on a manifold. These frames are assumed to be represented by vertices on a path graph  $P_n$ , where  $n$  is the frame number. The relationship between adjacent frames is modeled by the adjacency matrix  $\mathbf{W} \in \{0, 1\}^{n \times n}$  with  $W_{i,j} = 1$  if they are two consecutive frames,  $\|i - j\| = 1$ , and 0 otherwise. The

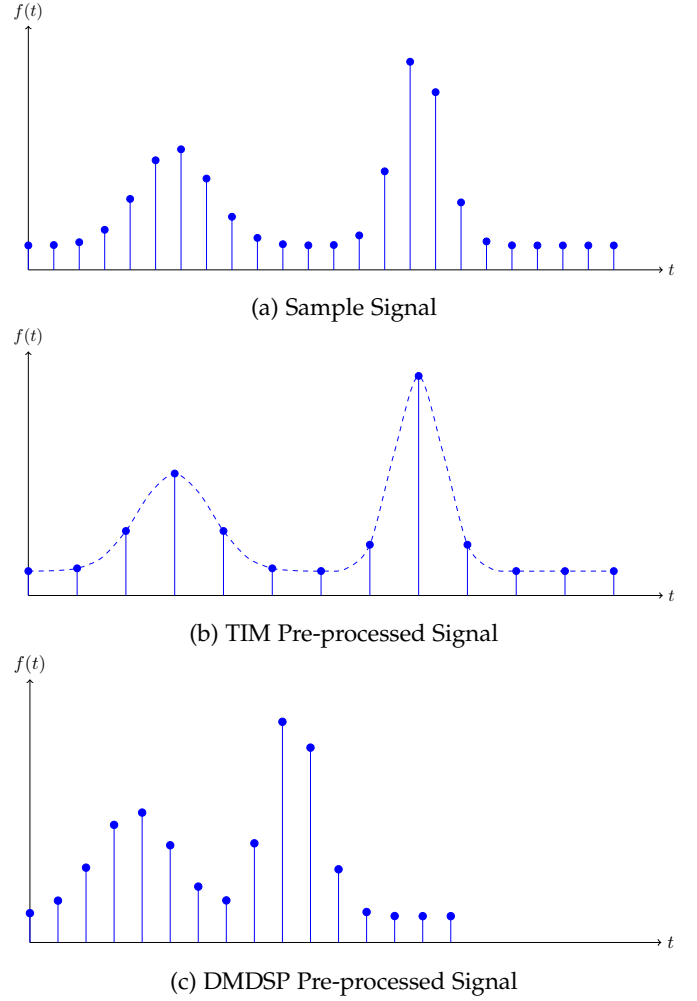


Fig. 3: Demonstration of how a sample signal (a) is pre-processed by TIM (b) and DMDS (c) where the vertical axis,  $f(t)$ , represents the magnitude of dynamics at the horizontal axis  $t$ .

graph  $P_n$  lies on a manifold when the total length of edges between connected vertices is minimized according to the following equation,

(1):

$$\sum_{i,j} (y_i - y_j)^2 W_{i,j}, \quad i, j = 1, 2, \dots, n \quad (1)$$

where  $\mathbf{y} = (y_1, y_2, \dots, y_n)$  is an eigenvector of the Laplacian graph of  $P_n$ . All points on an eigenvector  $\mathbf{y}_k$  are assumed to stay on a sinusoidal graph, formulated by,

$$f_k^n(t) = \sin(\pi kt + \pi(n - k)/(2n)), t \in [1/n, 1] \quad (2)$$

For  $n$  vertices, the Laplacian graph has  $n - 1$  eigenvectors  $\{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_{n-1}\}$ . In general, the manifold can be characterized with a collection  $F^n(t)$  of  $n - 1$  sinusoidal curves; moreover, frames at arbitrary positions are able to be interpolated as well. As each vertex in the graph  $P_n$  corresponds to a frame of an image sequence, the specific  $F^n$  of that sequence can be parametrized by mapping each frame to each point, defined by  $F^n(1/n), F^n(2/n), \dots, F^n(1)$ . Moreover, the parameterization requires linear extension of graph

embedding [28] in which a transformation  $\mathbf{w}$  is learned to minimize the Equation 1 as follows:

$$\arg \min_w \sum_{i,j} (\mathbf{w}^T \mathbf{x}_i - \mathbf{w}^T \mathbf{x}_j)^2 W_{i,j} \quad (3)$$

where  $x_i = \varepsilon_i - \bar{\varepsilon}$  is a mean-removed vector and  $\varepsilon_i$  is a vectorized image. He et al. [29] reformulated the minimization as an eigenvalue problem and solved it by singular value decomposition. More details about this minimization can be found in [27]. As mapping with frames of an image sequence defines a specific set of curves  $F^n(t)$ , synthesis of new frames at arbitrary temporal points is shown by Zhou et al. [27]. The synthesized frames seem to be a temporally smoothed version of the original sequence [6].

## 2.2 Sparse Sampling Approach: Sparsity-Promoting Dynamic Mode Decomposition (DMDSP)

The uniform selection approach, described in the previous Subsection 2.1, only temporally and regularly down-samples image sequences but does not eliminate redundant dynamics. TIM uniformly samples along a low-dimensional dense manifold, which is assumed to be the representative space of frames in a video sequence. The assumption would unknowingly eliminate some sparse structures of the dynamics. Therefore, alternative techniques for dynamic analysis could be used instead of graph-embedding manifold. A *sparse selection* approach tries to learn sparse structures of underlying dynamics and their appropriate magnitudes. Dynamic Mode Decomposition (DMD), which was first designed for capturing the momentum of indefinite-dimensional systems such as fluid flow, projects the complex system onto a low-complexity subspace spanned by dynamic modes with a few degrees of freedom. DMD does not make any rigid assumption about the existence of a manifold governing the overall dynamics, but freedom of an achieved model is data-driven. As such, DMD is only suitable for analyzing temporal dynamics but not learning their sparse structures. Therefore, we use a sparsity-constrained variant of DMD, i.e. the Sparsity Promoting DMD (DMDSP), developed by Jovanovic et al. [8] to select only significant dynamics of subtle expressions as visually illustrated in Figure 2. Further details of DMD and DMDSP are described in the following Sub-sections 2.2.1, 2.2.2.

### 2.2.1 Dynamic Mode Decomposition

Dynamic Mode Decomposition (DMD) [7] is designed to extract coherent structures at a single temporal frequency or dynamic mode e.g. flows in fluid dynamics [7] and motions in surveillance videos [21]. The DMD technique was first proposed and utilized in the analysis of fluid dynamics imagery. It analyzes sequences of “snapshots” from a data matrix, which are regularly sampled from fluid motions across time. Lets denote the  $N + 1$  sequential frames as  $\{\psi_0, \psi_1, \dots, \psi_N\}$ ; the previous frame i.e.  $\psi_0$  evolves to the next frame i.e.  $\psi_1$  over a regular temporal grid with a constant duration  $\Delta t$ . To model that evolution across all frames of an image sequence, two clusters of previous frames and next frames are formed as follows.

$$\Psi_0 := [\psi_0, \psi_1, \dots, \psi_{N-1}] \quad \Psi_1 := [\psi_1, \psi_2, \dots, \psi_N]$$

DMD assumes that dynamics between consecutive frames are governed by a linear time-invariant transformation  $A$  such that  $\psi_{t+1} = A\psi_t$  for  $t \in [0, \dots, N-1]$ . Hence,

$$\begin{aligned} \Psi_1 &= [\psi_1, \psi_2, \dots, \psi_N] \\ &= [A\psi_0, A\psi_1, \dots, A\psi_{N-1}] = A\Psi_0 \end{aligned} \quad (4)$$

For a rank- $r$  matrix of the cluster  $\Psi_0$ , the transformation  $A$  can be further spanned on a proper orthogonal basis  $U$  for an optimal representation  $F \in C^{r \times r}$  as follows.

$$A \approx U F U^* \quad \Psi_0 = U \Sigma V^* \quad (5)$$

where  $U^*$  is a complex conjugate transpose of the basis  $U$ , which is obtained from an economy-size singular value decomposition (SVD) of  $\Psi_0$ . The economy-size SVD is applicable for a tall matrix,  $\Psi_0 \in C^{M \times N}$ , as the number of pixels  $M$  of each frame  $\psi_t$  is often many more than  $N$ , the number of frames. As DMD models the evolution of the cluster  $\Psi_0$  into the cluster  $\Psi_1$  with a linear transformation  $A$ , the evolution can be regarded as a time-invariant system  $\Psi_1 = A\Psi_0$ . The transformation  $A$  can be found by minimizing the following Frobenius norm:

$$\arg \min_A \|\Psi_1 - A\Psi_0\|_F^2 = \arg \min_F \|\Psi_1 - U F \Sigma V^*\|_F \quad (6)$$

With a few linear algebra calculations, the optimal solution for the above formula is achieved as:

$$F_{dmd} = U^* \Psi_1 V \Sigma^{-1} \quad (7)$$

As  $F_{dmd}$  is a rank- $r$  matrix, it has a full set of linearly independent eigenvectors  $\{\mathbf{y}_1, \dots, \mathbf{y}_r\}$  and corresponding eigenvalues  $\{\mu_1, \dots, \mu_r\}$ . Then,  $F_{dmd}$  can be expressed in diagonal form,

$$F_{dmd} = Y D_\mu Z^* = [\mathbf{y}_1 \quad \dots \quad \mathbf{y}_r] \begin{bmatrix} \mu_1 & & \\ & \ddots & \\ & & \mu_r \end{bmatrix} \begin{bmatrix} \mathbf{z}_1^* \\ \vdots \\ \mathbf{z}_r^* \end{bmatrix} \quad (8)$$

where  $\{\mathbf{z}_1, \dots, \mathbf{z}_r\}$  are eigenvectors of  $F_{dmd}^*$ . These eigenvectors are bi-orthogonal to the  $\{\mathbf{y}_1, \dots, \mathbf{y}_r\}$ , which means  $Z^* Y = I$  and  $F_{dmd}^t = Y D_\mu^t Z^*$ . As a linear time-invariant system  $A$  is assumed to govern the dynamics between sequential frames, an evolution of a frame  $\psi_t$  can be formulated according to an initial frame  $\psi_0$  as below.

$$\psi_t = A^t \psi_0 \approx (U F_{dmd} U^*)^t \psi_0 = U Y D_\mu^t Z^* U^* \psi_0 \quad (9)$$

Let  $\Phi = U Y$  be the DMD modes and  $A = Z^* U^* \psi_0$  as the corresponding magnitudes, then a frame  $\psi_t$  can be rewritten regardless of  $\psi_0$  as  $\psi_t = \Phi D_\mu^t A$ . With the amplitudes in  $D_\alpha$  and the Vandermonde matrix,  $V_{and}$ , representing temporal evolution and  $\Phi$  representing spatial modes, the frame cluster  $\Phi_0$  is reformulated as follows:

$$\begin{aligned} \Psi_0 &= [\psi_0 \dots \psi_{N-1}] \approx \Phi D_\alpha V_{and} \\ &= [\phi_0 \dots \phi_r] \begin{bmatrix} \alpha_1 & & \\ & \ddots & \\ & & \alpha_r \end{bmatrix} \begin{bmatrix} 1 & \dots & \mu_1^{N-1} \\ \vdots & \ddots & \vdots \\ 1 & & \mu_r^{N-1} \end{bmatrix} \end{aligned} \quad (10)$$

Determination of unknown amplitudes  $\alpha := [\alpha_1 \dots \alpha_r]^T$  depends on the solution of the following optimization problem.

$$\arg \min_{\alpha} \|J(\alpha)\| = \arg \min_{\alpha} \|\Psi_0 - \Phi D_{\alpha} V_{and}\|_F^2 \quad (11)$$

where the optimal DMD amplitudes can be obtained by,

$$\alpha_{dmd} = ((Y^* Y) \circ (V_{and} \bar{V}_{and}^*))^{-1} \text{diag}(V_{and} V \Sigma^* Y) \quad (12)$$

### 2.2.2 Sparsity Promoting Dynamic Mode Decomposition

Though an optimal amplitude for each mode is found by the Equation (11), the amplitudes are solved with the assumption that every mode is equally significant. However, it is not always true and especially not correct for dynamics of micro-expressions. These dynamics are sparse as these expressions are naturally concise and sudden. Therefore, only these sparse modes are significant while the rest can be removed without much loss in re-construction of the original signals. In order to reveal the sparsity of the dynamics and its appropriate amplitudes, we adopt the sparsity-promoting DMD approach of Jovanovic et al. [8] which adds sparse constraints into the DMD formulation. It allows trade-off between loss of signal reconstruction and the number of sparse modes.

As DMD only analyzes data and does not apply any sparse constraints, the number of DMD modes is equal to the number of frames. DMDSP selects a subset of these DMD modes which have dominant influence on the reconstruction of a given sequence. Implementation of sparsity-constraints involves the following two steps.

- 1) Identification of sparsity structure such that a user-defined trade-off between the number of extracted modes and approximation error is achieved.
- 2) Identification of optimal amplitudes for extracted modes given the sparsity structure.

Jovanovic et al. [8] suggests that the sparse structure problem in the first step can be relaxed and formulated as the  $l_1$ -norm of the vector of magnitudes  $\alpha$ :

$$\arg \min_{\alpha} J(\alpha) + \gamma \sum_{i=1}^r |\alpha_i| \quad (13)$$

where  $\gamma$  is a sparsity regularization parameter and  $\alpha_i \in \alpha$  is a DMD magnitude at rank- $i$ , showing sparseness of the vector  $\alpha$ . The Alternative Direction Method of Multipliers (ADMM) method [30] is utilized to solve the above convex optimization problem. Given the fixed sparsity structure, vectors of magnitudes  $\alpha$  can be optimized as a solution to the following constrained convex optimization problem:

$$\arg \min_{\alpha} (|J(\alpha)|) \quad \text{s.t.} \quad E^T \alpha = 0 \quad (14)$$

where the matrix  $E$  represents the sparsity structure of the amplitude vector  $\alpha \in R^{r \times m}$ , identified in the first step;  $m$  represents the number of  $\alpha_i$  with zero values. Each column vector has only one non-zero element corresponding to each zero component of  $\alpha$ ; for example, with  $\alpha \in \mathbb{C}^4$  and  $\alpha = [\alpha_1, 0, \alpha_3, 0]^T$ ,  $E$  is given as follows.

$$E^T = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (15)$$

Jovanovic et al. [8] show that the optimal DMD amplitudes with a fixed sparsity structure can be computed as follows.

$$\alpha_{dmdsp} = [I \quad 0] \begin{bmatrix} P & E \\ E^T & 0 \end{bmatrix}^{-1} \begin{bmatrix} q \\ 0 \end{bmatrix} \quad (16)$$

More details about solutions of these sparsity-constrained problems can be found in [8].

### 2.3 Sparse & Uniform Sampling for Redundancy Removal

In the uniform approach, TIM partially removes redundant information by interpolating or synthesizing fewer frames at an arbitrarily temporal grid. This new grid is defined by the number of interpolated frames such that they are temporally equispaced. This technique of redundancy removal only works if more frames are interpolated at the significant part and less frames are interpolated at the insignificant part of a sample signal. For example, the discrete signal of the Figure 3a has 24 points nine of which are in a significant region and the rest are not. When the uniform sampling interpolates the signal and synthesizes only 13 points in which six are in significant regions and the rest are not. The ratio between redundant and significant points is dropped from  $\frac{15}{9}$  of the original signal to  $\frac{7}{6}$  of the synthesized signal. It means that the generated signal becomes less redundant. The significant drawback of this approach is that it interpolates new frames at regularly separated positions regardless of the signal's sparse structure. As positions of interpolated frames on the signal are decided by how many frames should be interpolated, that number can only be arbitrarily decided without any knowledge about structures of input signals. Moreover, even with this prior knowledge, finding optimal numbers of generated frames is also difficult due to regularly spaced temporal interpolation. The interpolation is unadaptive to the signals' structures. In brief, there is no specific method to guarantee optimal removals of redundancy in the case of uniform sampling.

The sparse sampling approach, DMDSP, tackles significant drawbacks in the uniform sampling approach, TIM. While DMD, like TIM, only analyzes the signal and models dynamics in a sample signal regardless of its sparse structure, DMDSP (sparse-promoting DMD) has incorporated the sparse constraints into the analysis. It is formulated as a convex optimization problem in Equations (14) and (16). Solutions of these problems are sparse structures  $E$  of the sample signal and their optimal amplitudes  $\alpha$ . With the sparsity constraint, amplitudes of dynamics are large if they have profound contribution on the approximate reconstruction of original sequences. Otherwise, their amplitudes are small or nearly zero. In the sparse sampling approach, redundancy is eliminated by removing modes with small or near-zero amplitudes; then, sequences are reconstructed with the remaining dominant modes. The reassembled signal is shorter and only contains significant parts of the signal, as demonstrated in Figure 3c, a result of applying sparse sampling on the sample signal of Figure 3a. The reconstructed signal in Figure 3c demonstrates the main advantage of sparse sampling over the uniform one as redundancy of the dynamics in a signal is removed more accurately by DMDSP than TIM.



### 3 DYNAMIC ANALYSIS OF SPONTANEOUS SUBTLE EMOTIONS

Since human beings control several muscles, either consciously or sub-consciously, beneath their skins to form facial expressions as the main non-verbal communication of their emotional states, characteristics of the expressions heavily depend on muscular movements. Understandably, researches of automatic facial expression recognition show superiority of spatiotemporal features over their spatial counterparts i.e. LBPTOP has been shown to be better than LBP and Gabor-like features in [22]. It shows that temporal changes or dynamics across consecutive frames significantly contribute to recognition of normal facial expressions. For subtle expressions, careful analysis of temporal dynamics becomes more crucial as subtle facial expressions are not visually far from neutral faces, especially to human eyes. The possibility of correctly recognizing or detecting subtle emotions greatly depends on how discriminating the temporal dynamics across consecutive frames are. Despite their importance, a spatiotemporal feature like LBPTOP only partly utilizes these dynamics without further analysis, visualization or enhancement. To better leverage these subtle facial movements, we utilise the DMD analysis technique, introduced in 2.2.1, to separate, analyze and visualize facial dynamics. As TIM and DMDSP are only sampling the original sequence uniformly or sparsely, respectively, the result of these pre-processes are sequences of generated frames in either case. Therefore, like original frames, each generated frame can be then analyzed by DMD for their contributions to overall dynamics regardless of sampling strategies. DMD decomposes an image sequence into three main components: spatial modes  $\Phi$ , amplitudes  $\alpha$  and temporal dynamics  $\mu$  as per Equation (10). While spatial modes and temporal dynamics represent separated spatial and temporal information, an amplitude measures a contribution of each spatial mode and its corresponding temporal dynamic of each frame in the reconstruction of the original image sequence. Since frames are assumed to be linearly independent, the number of dynamic modes and magnitudes, i.e. the rank  $r$ , is equal to the number of input frames. As a result, each frame has its corresponding DMD amplitude. The more activities each frame of facial expressions has, the larger its DMD magnitude is. DMD magnitude is a frame-by-frame measurement of facial activities; therefore, we can plot the amplitudes against the frame index as a visualization of micro-expressions' dynamics. These plots show temporal analyses of dynamics for video sequences of subtle emotions from the CASME II and SMIC corpora. A disadvantage of the above visualization method is that this temporal analysis is only applicable for individual video sequences. However, it is not able to be generalized for the whole databases collectively due to the variety of frame-lengths among video sequences. To visualize the general dynamics of subtle emotions in a database, we propose spectral analyses in place of temporal analyses as the frequency bandwidth of dynamics is fixed to a certain range with respect to sampling rates or recording frame rates regardless of frame-lengths. With temporal and spectral analyses, we can visualize and analyze the dynamics of both individual samples and whole databases. When redundancy of dynamics is gradually re-

moved, these analyses provide observations of changes in dynamics at both levels. Therefore, they help identify which temporal parts and frequency bandwidths are important or redundant for a video sequence of a micro-expression or a whole database of spontaneous subtle emotions. These observations of significant temporal and spectral parts of signals can be compared with prior knowledge about micro-expressions such as the duration of micro-expressions lasting between  $\frac{1}{25}$ s and  $\frac{1}{15}$ s [5]. More importantly, they provide visualization of dynamic behaviors for several classes of subtle expressions. Hence, we can learn about temporal and spectral dynamical characteristics of each type of subtle emotion across different levels of redundancy removals or dynamic sparsity.

#### 3.1 Temporal Analysis

Analysis of dynamics in the temporal domain involves plotting the DMD amplitudes of an image sequence against corresponding frame indexes. The plots of amplitudes over time give a hint at temporal locations where most motions happen. In this section, the temporal analyses are done for sequences, pre-processed by uniform and sparse sampling approaches of redundancy removals for several percentages of preserved frames. They help analyze and visualize how sparse and uniform sampling affect overall dynamics of sequences. The following subsections demonstrate without loss of generality the temporal analysis of uniform and sparse dynamics for two randomly selected samples from the CASME II database ( EP02\_01f of Subject 01, EP08\_03 of Subject 17 ).

##### 3.1.1 Temporal Analysis of Uniform Sampling

In the uniform sampling, dynamics are learned through interpolation of an original sequence and regenerated into a shorter one. This generated sequence is then analyzed by DMD for the dynamic amplitudes of corresponding frames. As TIM regularly samples an original dynamic by a temporally equispaced grid, the DMD amplitudes can be mapped to the original sequence in the same grid. For example, lets assume an input sequence of 32 frames; TIM interpolates five frames at the following indexes: 1,8,16,24,32. Then DMD amplitudes of these five frames correspond to dynamics at five positions 1,8,16,24,32 of the original sequence while other positions are supposed to have zero-dynamics. The number of interpolated frames over the original frame length is a percentage of preservation. The above approach is used for computing and plotting temporal dynamics of the chosen sequences in Figures 4, 5. In these plots, 25%, 50%, 75%, 95% and 100% of the frames are preserved for temporal analysis. Note that 100% means no TIM interpolation of original frames has been done and a video sequence is directly analyzed by DMD for dynamical amplitudes. As TIM interpolation likely down-samples the dynamics, the temporal dynamics of lower percentages of preserved frames look like the blurred versions of higher percentages.

##### 3.1.2 Temporal Analysis of Sparse Sampling

As the sparse sampling approach looks for significant frames of original dynamics, it is possible to determine

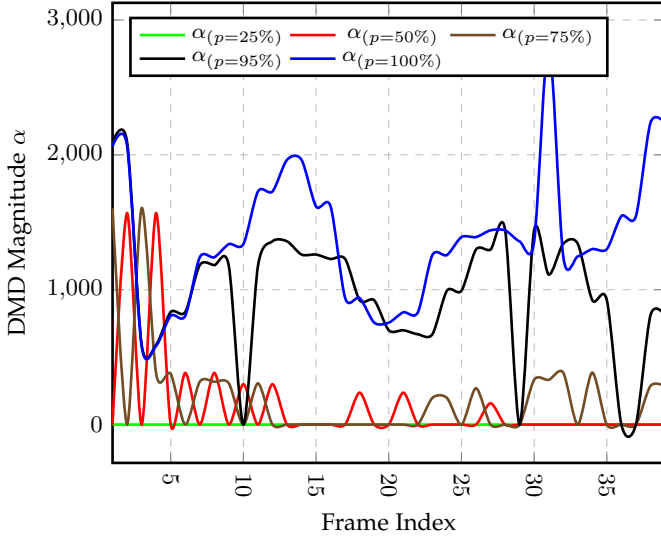


Fig. 4: Temporal Analysis with Uniform Sampling in which plots illustrate facial dynamics of Subject 01 in the sample EP02\_01f of the CASME II corpus w.r.t. percentages ( $p$ ) of preserved frames: 25%, 50%, 75%, 95% and 100% of the original frame length

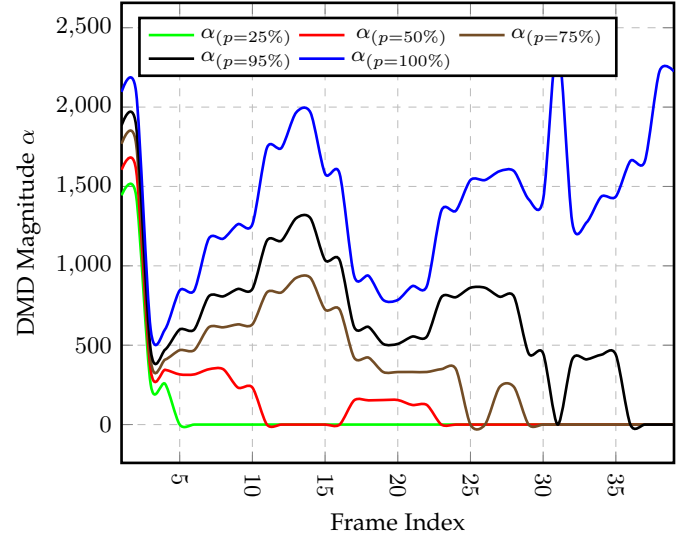


Fig. 6: Temporal Analysis with Sparse Sampling in which plots illustrate facial dynamics of Subject 01 in the sample EP02\_01f of the CASME II corpus w.r.t. percentages of preserved frames: 25%, 50%, 75%, 95% and 100% of the original frame length

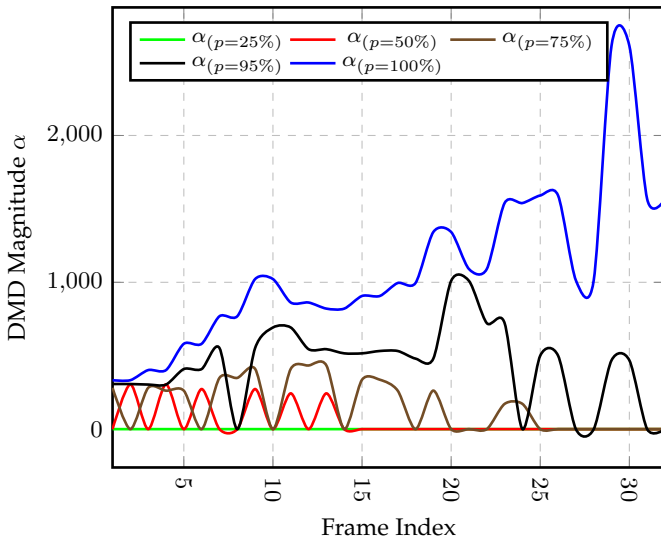


Fig. 5: Temporal Analysis with Uniform Sampling in which plots illustrate facial dynamics of Subject 17 in the sample EP08\_03 of the CASME II corpus w.r.t. percentages ( $p$ ) of preserved frames: 25%, 50%, 75%, 95% and 100% of the original frame length

redundant frames among original sequences. Correspondingly, amplitudes of these frames are masked to zeros while the rest of the amplitudes are analyzed in DMD. Since sparse structures are regulated by the  $\gamma$  parameter, its value also controls how many percentages of frames are preserved. For a certain range of  $\gamma$  value [38 – 20000], these percentages vary between 5% and 100%. Figures 6 and 7 show temporal analyses of sparsely sampled dynamics for 25%, 50%, 75%, 95% and 100% preserved frames. Note that 100% means the DMD amplitudes of original sequences are not processed by DMDSP. These plots show how the sparse dynamics change

across different percentages. Furthermore, they show consistency of sparse structures, learned by DMDSP with different  $\gamma$  values, especially in Figure 7.

Comparisons between temporal analyses of uniform and sparse sampling (Figure 4 vs Figure 6 or Figure 5 vs Figure 7) show that sparse dynamics have more consistent and clearer patterns than uniform dynamics. For instance, plots of DMD magnitudes of the sample EP08\_03 have consistent shapes for frames 0-19 across multiple percentages (25%, 50%, 75%) of preservation. Meanwhile, uniform sampling blindly selects equispaced indexes and removes in-between frames. Therefore, plots of DMD magnitudes  $\alpha$  are radically different given various percentages of preservation  $p$  in the Figure 5. Furthermore, the DMD magnitudes  $\alpha$  in the Figure 5 represent how much each frame contributes into the reconstruction of original dynamics. Therefore, consistent structures of sparse sampling plots, Fig. 7, and their magnitudes provide locations and weights of significant dynamics in an image sequence.

### 3.2 Spectral Analysis

Temporal analysis is only useful for inspecting dynamics of a single sample. It is not able to summarize and visualize dynamics of several samples due to varying frame-lengths across different samples. Therefore, temporal analysis is unable to show dynamical characteristics of the whole CASME II or SMIC databases. Spectral analysis avoids the frame-length problem through analyzing video samples of databases in the frequency domain. As DMD decomposes an image sequence into spatial modes  $\phi$ , amplitudes  $\alpha$  and temporal dynamics  $\mu$  in Equation (10), the frequency of dynamics  $f_{DMD}$  can be computed from the temporal dynamics  $\mu$  as follows.

$$f_{DMD} = \log(\mu) f_s \quad (17)$$



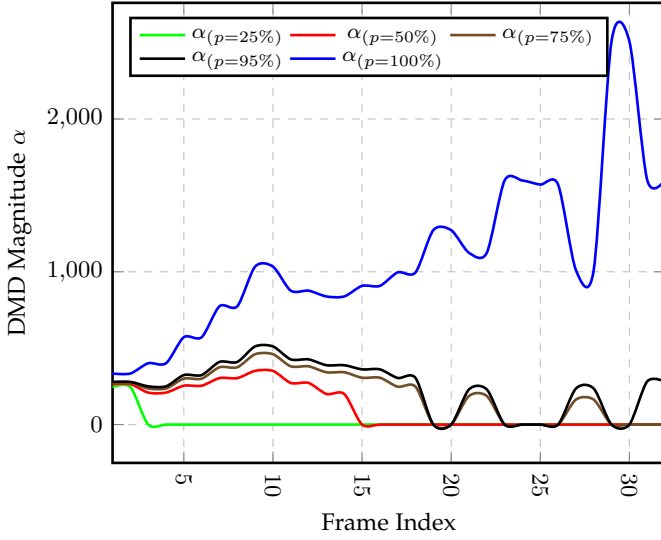


Fig. 7: Temporal Analysis with Sparse Sampling in which plots illustrate facial dynamics of Subject 17 in the sample EP08\_03 of the CASME II corpus w.r.t percentages of preserved frames: 25%, 50%, 75%, 95% and 100% of the original frame length

where  $f_s$  is the maximum sampling frequency or frame-rate of an image sequence. The frame-rate also limits the bandwidth of dynamics as  $f_{DMD} \in B_{DMD} = [0 \dots f_s/2]$  Hz according to the Nyquist-Shannon Sampling Theorem [2] which states that the maximum frequency of a signal is half of its sampling frequency. For example, the CASME II database is recorded with frame-rate 200 fps; then, its spectral bandwidth of dynamics is  $[0 - 100]$  Hz. Similarly, SMIC has  $[0 - 50]$  Hz bandwidth according to its frame-rate 100 fps. With the Equation (17), a frequency value  $f_{DMD}$  can be found for each temporal dynamic value  $\mu_i$ ; therefore, each  $f_{DMD}$  has a corresponding amplitude  $\alpha$ . Moreover, spectral analysis of an image sequence results in a histogram of amplitudes  $\alpha$  over bins of frequencies  $f_{DMD}$ . Given that all videos in the same database have the same frame-rate, bandwidth, hence the same number of bins, it is valid to sum up histograms of all sequences from the same database to produce spectral analysis of a database. For instance, Figure 8 shows the spectral analysis with DMD of CASME II and Figure 9 demonstrates the dynamics of SMIC over its spectral bandwidth. These figures show the spectra of temporal dynamics from original videos of CASME II and SMIC, which are not pre-processed by neither uniform nor sparse sampling. These spectra of both databases are spread over the whole bandwidth; therefore they do not provide any information about dominant temporal dynamics of micro-expressions. The lack of dominant spectra might be due to noises, generated by redundancies of neutral faces as their differences are not due to motions but illuminating conditions.

The proposed sparse sampling techniques remove redundant neutral frames and reveal significant spectra of dynamics. In addition, plots of spectral analysis visualize the dominant spectra and how a whole database responds to different sparsity levels or percentages of preserved frames.

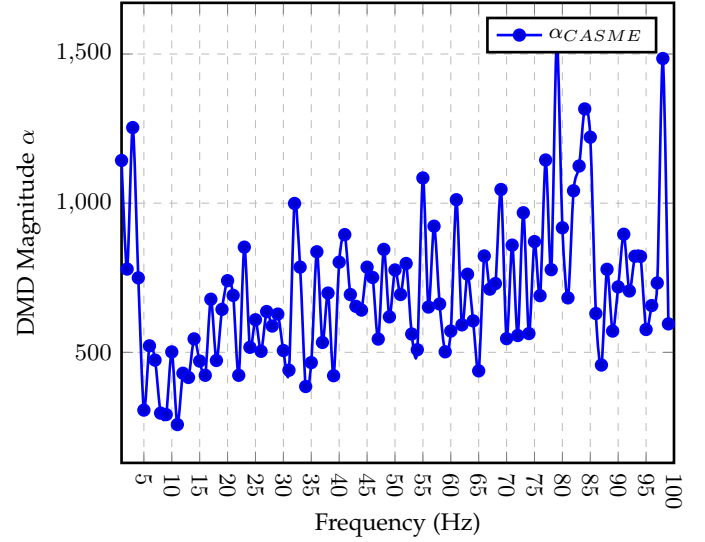


Fig. 8: Spectral Analysis on Spontaneous Subtle Emotions of CASME II

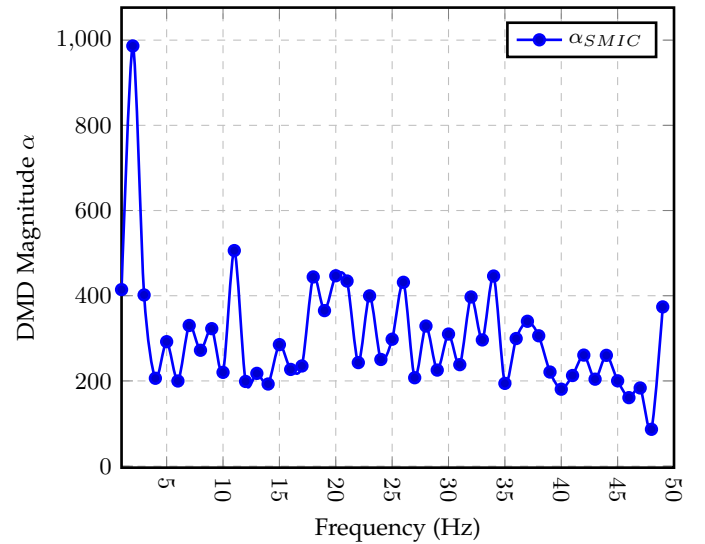


Fig. 9: Spectral Analysis on Spontaneous Subtle Emotions of SMIC

In the following subsections, spectral analyses are used for analyzing uniformly and sparsely sampling approaches with respect to different percentages. Note that these percentages are inversely proportional to the sparsity levels. Five percentage values i.e. 5%, 25%, 50%, 75% and 100% of preserved frames over the original frame-lengths are used in the following analysis of both uniform and sparse sampling.

### 3.2.1 Spectral Analysis of Uniformly Sampling Dynamics

Section 2.3 discussed the technicality of the uniform sampling approach like TIM [6] for removing dynamic redundancies. In this section, we visualize redundancy removals of the uniform sampling in the spectral domain for whole databases. Figures 10 and 11 demonstrate several spectral analyses of CASME II and SMIC with respect to four percentages of preserved frames  $p = \{5\%, 25\%, 50\%, 100\%\}$ , four plots with different colors. These plots in Figures 10

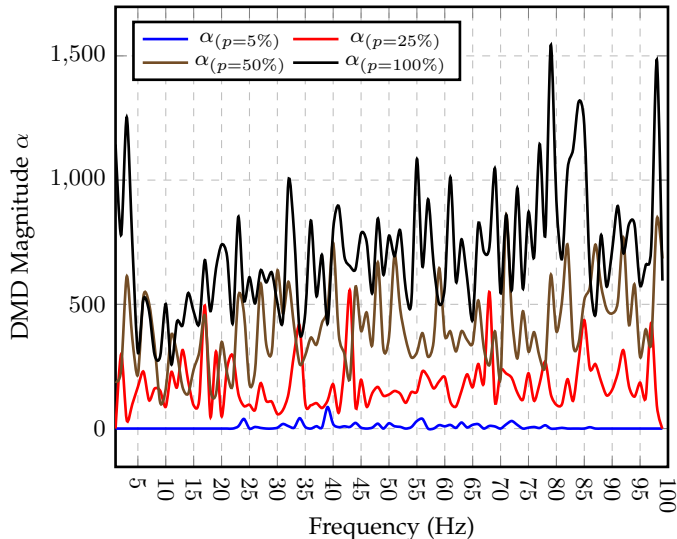


Fig. 10: Spectral Analysis of CASME II database pre-processed with Uniformly Sampling approach and four percentages of preserved frames  $p = \{5\%, 25\%, 50\%, 100\%\}$

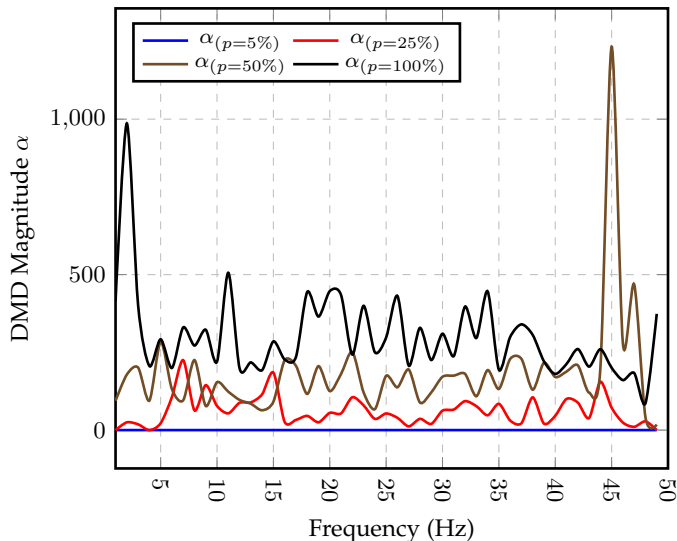


Fig. 11: Spectral Analysis of SMIC database pre-processed by Uniformly Sampling approach and four percentages of preserved frames  $p = \{5\%, 25\%, 50\%, 100\%\}$

and 11 show that the uniform sampling approach tends to suppress the dynamics of the lowest and highest parts of the bandwidth while they smoothen the dynamics in the middle frequency ranges. In other words, the uniform sampling approach behaves like a band-pass filter for dynamics since frames are interpolated at equispaced points along a manifold. It linearly combines nearby original frames, preserves spatial characteristics whilst smoothing temporal profiles [6] [27].

### 3.2.2 Spectral Analysis of Sparsely Sampling Dynamics

In contrast to the previous section 3.2.1 about spectral analyses for uniform dynamics by TIM [6] on CASME II and SMIC databases, this section shows spectral analysis of sparse dynamics with DMDSP [8]. While Section 2.2.2

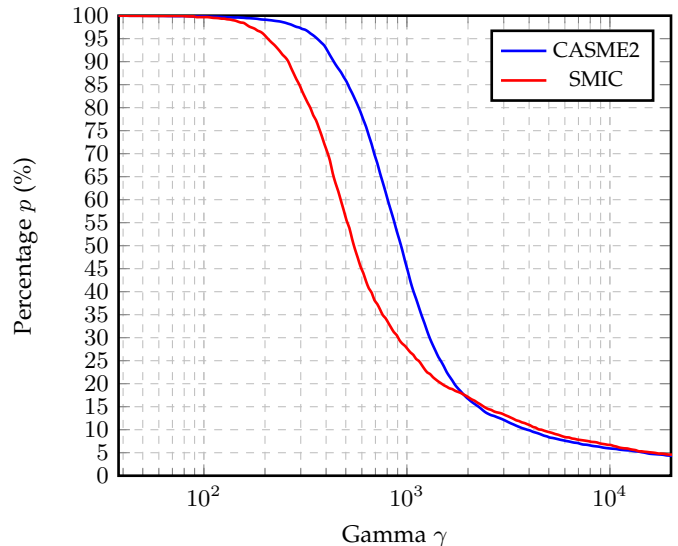


Fig. 12: Gamma Analysis of CASME II and SMIC databases with Sparse Sampling approach, which demonstrates inversely proportional relation between percentages of preserved frames  $p$  and sparsity parameter  $\gamma$

elaborates how the sparsity constraint with a regulation parameter  $\gamma$  is formulated in DMDSP, this section explains a relation between sparsity levels and percentages of preserved frames. More importantly, it analyses the sparse dynamics on CASME II and SMIC databases.

The regulation parameter  $\gamma$  in Equation (13) controls sparseness enforced by DMDSP. The higher the  $\gamma$  values are, the more sparsity constraints are enforced. In other words, more dynamics are removed with increments of  $\gamma$  values. However, the optimization problem in Equation (13) also involves losses during reconstruction from the sparse dynamics. To achieve this optimization, trade-offs are made between reconstruction loss and the number of dynamics. The solution is to increasingly removing unimportant frames from the sequence as the sparsity level i.e.  $\gamma$  increases. Hence, there is an inversely proportional relation between sparsity values  $\gamma$  and percentages of preserved frames. This relationship is shown through plots of percentages over logarithmic ranges of  $\gamma$  values in Figure 12. To generate the Figure 12, we choose 400 gamma values, equispaced in the logarithmic scale, in a range [38-20000]. DMDSP is applied to a video sample with these gamma values. Due to the sparsity constraint, many suppressed dynamics with near-zero amplitudes are removed from the sequences. From the number of remaining frames, a percentage of preserved frames can be computed for the corresponding gamma value. Figures 13 and 14 show spectral analyses of CASME II and SMIC databases with different percentages of preserved frames for redundancy removals. Spectral analyses are done for 5%, 10%, 15% until 100%; however, only a few are displayed for simplicity of the plot. Sparse structures are learned with DMDSP, sparsity-promoting DMD, on dynamics of image sequences in the sparse sampling while the uniform sampling with TIM just averagely interpolates from original frames. Spectral analyses of both databases show similar trends in which high-frequency components

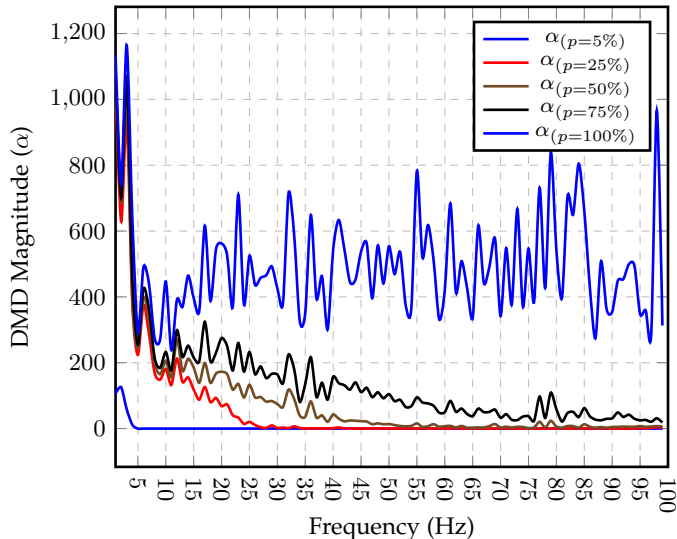


Fig. 13: Spectral Analysis of CASME II database, pre-processed by Sparse Sampling approach and five percentages of preserved frames  $p = \{5\%, 25\%, 50\%, 75\%, 100\}$

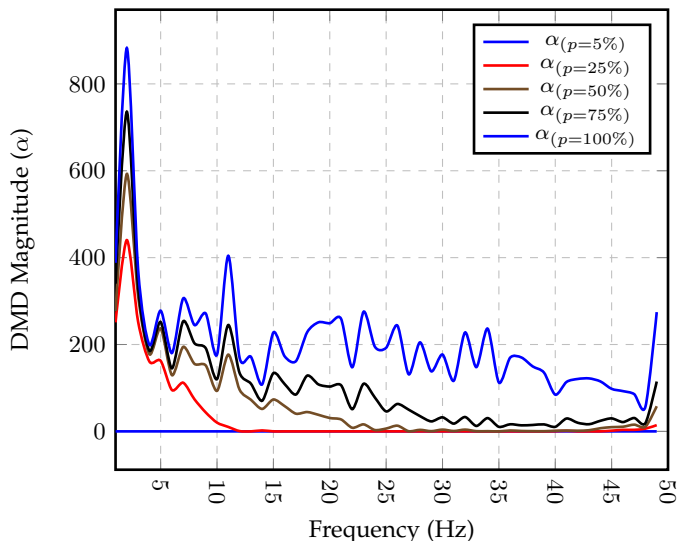


Fig. 14: Spectral Analysis of SMIC database, pre-processed by Sparse Sampling approach and five percentages of preserved frames  $p = \{5\%, 25\%, 50\%, 75\%, 100\}$

are suppressed by the decreasing percentages or increments of sparsity. In other words, most sparse structures or dynamically significant frames are located at low frequencies from 0 Hz to 25 Hz. According to Ekman [5], the duration of subtle emotions is more than  $\frac{1}{25}$  s. Therefore, the spectra of temporal dynamics should be significant below 25 Hz, which are well spotted by the sparse sampling approach in both CASME II and SMIC databases.

## 4 EXPERIMENTS & DISCUSSION

In this section, our experimental results for different approaches of redundancy removal are reported and analyzed for both the CASME II and SMIC databases. The CASME II corpus [3] has samples from seven different classes or labels:

“others” (O), “disgust” (D), “happiness” (H), “surprise” (S), “repression” (R), “fear” (F), “sadness” (S). However, only the first five of the mentioned classes have enough samples for any statistically meaningful experiments, as shown in the Table 2, and as considered in the database designers’ baseline experiments. Meanwhile, the SMIC corpus includes three different emotional categories: “positive” (P), “negative” (N) and “surprise” (S). Hence, we report experimental results of 5-class and 3-class classification for CASME II and SMIC corpora accordingly. For both CASME II and SMIC, labeling was done by the designers per a video sample rather than per a frame. In other words, all frames of a video sample are assumed to have the same emotional labels as each video sample is actually a cropped video clip from a raw footage based on detected onset (beginning of the emotion) and offset (ending) points. Hence, the automatic recognition of subtle emotions is trained and tested at the granularity of video samples assuming that facial expressions in a video sample correspond to only one emotional state. The assumption is justified by procedures of sample acquisition in CASME II [3] and SMIC [6].

In the following experiments, “Sparse Sampling” (SS), “Uniform Sampling” (US), “Generative” (US\*) and “Random” (RA) approaches for redundancy removals are compared against the “Baseline” (BL) result of subtle emotion recognition tasks. Note that the US and US\* approaches both utilize TIM for generating temporal dynamics but for different frame-length parameters. In the US approach, the number of frames is adaptively synthesized with respect to percentages (45%, 50%, ..., 100%) of preserved frames for each sequence. The US\* approach strictly interpolates each video sample into fixed frame-lengths i.e. 150 frames in Wang et al. [16] and 10 frames in Pfister et al. [6]. Meanwhile, the BL approach shows performances of systems without eliminating dynamic redundancies. In the RA approach, frames are randomly collected from the original sequences such that a certain percentage (45%, 50%, ..., 100%) of frame-length is acquired. While the SS and US approaches are discussed in detail and comparatively analyzed in previous sections 2.1 and 2.2, inclusion of the RA, BL, US and US\* approaches aims to confirm the superiority of the proposed Sparse Sampling (SS) approach over random or no elimination of redundant frames as well as uniform sampling approaches from the previous works [3], [6]. Experimental parameters are further elaborated for the rest of this section and summarized in Table 4.

TABLE 4: Experimental parameters of sampling approaches, uniform LBPTOP feature, and Support Vector Machine (SVM) classifier with respect to SS, US, US\*, RA and BL approaches

	Preprocessing		Uniform LBPTOP		SVM		
	Method	Parameter	CASME II	SMIC	K	c	g
SS	DMDSP	[45:5:100] %	5x5x1	8x8x1	RBF	$10^4$	.5
US	TIM	[45:5:100] %	5x5x1	8x8x1	RBF	$10^4$	.5
US*	TIM	CASME II: 150 SMIC: 10	5x5x1	8x8x1	RBF	$10^4$	.5
RA	RAN	[45:5:100] %	5x5x1	8x8x1	RBF	$10^4$	.5
BL	N/A	N/A	5x5x1	8x8x1	RBF	$10^4$	.5

TABLE 1: Performance measures (F1 score (F1), Recall Rate (R) and Precision Rate (P)) of the proposed **SS** approach, in comparison to the **US**, **US\***, (**RA**) and **BL** approaches w.r.t each class of CASME II (O,D,H,S,R) and SMIC (P,N,S) corpora

	CASME II												SMIC											
	Others (O)			Disgust (D)			Happiness (H)			Surprise (S)			Repression (R)			Negative (N)			Positive (P)			Surprise (S)		
	F1	RR	PR	F1	RR	PR	F1	RR	PR	F1	RR	PR	F1	RR	PR	F1	RR	PR	F1	RR	PR	F1	RR	PR
<b>SS</b>	.56	.64	.50	.27	.23	.32	.58	.55	.62	.67	.52	.93	.39	.41	.38	.53	.48	.59	.60	.70	.53	.64	.62	.67
<b>US</b>	.52	.58	.47	.09	.07	.12	.36	.48	.29	.34	.24	.60	.32	.30	.35	.37	.36	.38	.40	.46	.35	.44	.37	.55
<b>US*</b>	.47	.52	.44	.20	.19	.20	.37	.44	.32	.12	.08	.22	.18	.15	.24	.46	.41	.52	.49	.53	.45	.48	.51	.46
<b>RA</b>	.49	.54	.44	.13	.12	.14	.25	.42	.18	.10	.06	.19	.32	.30	.34	.38	.36	.41	.36	.46	.30	.37	.29	.52
<b>BL</b>	.47	.53	.42	.25	.22	.27	.33	.31	.34	.42	.40	.43	.30	.26	.35	.39	.36	.42	.41	.49	.35	.39	.35	.45

TABLE 2: Number of video samples for each class of CASME2 (H,D,R,S,O) and SMIC (P,N,S) corpora

CASME II			SMIC		
Emotion	Label	# samples	Emotion	Label	# samples
Happiness	H	33	Positive	P	51
Disgust	D	60	Negative	N	70
Repression	R	25	Surprise	S	43
Surprise	S	27			
Others	O	102			
Fear	N/A	2			
Sadness	N/A	7			

After dynamics of signals are manipulated and shorter videos are reconstructed from remaining dynamic modes, all video frames of the same database are spatially normalized as they are resized to a fixed resolution. Denote that the resized and shortened videos have  $n_r$  rows,  $n_c$  columns and  $n_f$  frames. In the CASME II corpus,  $n_r$  and  $n_c$  are set to 340 and 280 respectively; meanwhile,  $n_r$  and  $n_c$  for SMIC are fixed at 170 and 140 respectively. Then, a feature vector is extracted from each video sample as follows: (i) dividing each image sequence into a  $n_b \times n_b \times 1$  matrix of non-overlapping volumetric 3-D blocks with size  $(\frac{n_r}{n_b}, \frac{n_c}{n_b}, \frac{n_f}{1})$  where ( $n_b = 5$ ) for CASME II and  $n_b = 8$  for SMIC); (ii) extracting a histogram from each block by uniform 4-connected LBPTOP<sub>4, 4, 4, 1, 1, 3</sub> [22] with (1, 1, 3) as correspondent horizontal ( $R_x$ ), vertical ( $R_y$ ) and temporal ( $R_t$ ) radii; (iii) concatenating all histograms of  $n_b \times n_b$  volumetric cells into the feature vector. The number of blocks and LBPTOP configuration are chosen for replicating experimental results from Table 4 in [3] and Table 4 in SMIC [9]. Though Yan et al. [3] proposes that  $R_t = 4$  gives the best accuracy, we choose  $R_t = 3$  to satisfy minimum frame-length ( $2 \times R_t + 1 = 7$ ) for LBPTOP feature extraction even if only 45% of original frames are preserved. Given the mentioned details of LBPTOP feature extraction for CASME II and SMIC, dimensions ( $D \times 1$ ) of feature vectors are computed by  $D = n_b^2 \times 15 \times 3$ . Note that 15 is the dimension of each uniform LBPTOP histogram, 3 is the number of considered planes  $XY, XT$  and  $YT$ , and  $n_b^2$  is the number of blocks. As a result, the dimensions ( $D$ ) of features for CASME II and SMIC experiments are 1125 and 2880.

After LBPTOP feature extraction, Support Vector Machine (SVM) with Radial Basis Function (RBF) is utilized for training classifiers with extracted feature vectors. The RBF kernel is employed so that the classifier SVM optimizes decision hyper-planes for these multi-class recognition problems in the infinite similarity spaces of RBF kernel rather than original small feature spaces i.e. 1125-D feature of CASME

TABLE 3: Performance measures (F1 score, Recall Rate and Precision Rate) of the proposed **SS** approach, in comparison to the **US**, **RA** and **BL** approaches

	CASME II				SMIC			
	ACC	F1	RR	PR	ACC	F1	RR	PR
<b>SS</b>	.49	.51	.47	.55	.58	.60	.60	.60
<b>US</b>	.38	.35	.33	.37	.40	.41	.40	.43
<b>US*</b>	.33	.28	.27	.28	.48	.48	.49	.48
<b>RA</b>	.34	.29	.29	.29	.37	.39	.37	.41
<b>BL</b>	.38	.35	.34	.36	.40	.40	.40	.41

II samples and 2880-D features of SMIC samples. So far, the cost parameter ( $c$ ) and Gaussian Kernel Bandwidth ( $g$ ) are set to  $10^4$  and 0.5 respectively, as shown in Table 4. The kernel bandwidth is kept at a default value 0.5 and the cost parameter ( $c$ ) is set to a very large value to prevent misclassification of inter-class samples from the same subject.

As facial expressions and identities are inseparable characteristics of all video samples, the special Leave-One-Subject-Out (LOSO) protocol [31] [32] [9] [6] is also enforced to prevent subject identities interfere with classification of subtle emotions. LOSO is an exhaustive cross-validation technique, partitioning corpus testing and training sets according to subject identities. It requires samples of the testing corpus from one subject and those of the training corpus from the other subjects. For example, if Subject 1 of CASME II corpus is chosen as the test subject, all nine video samples of the subject are preserved for testing. The other video samples (238) are used for training the classifier SVM-RBF. These processes are folded until every subject is employed in the test corpus once. There are 26 subjects in the CASME II database hence 26 folds of cross-validation in the LOSO protocol are carried out. All 26 combinations of training and testing samples can be inferred from Table 2 in [4]. For performance metric, Accuracy (ACC), F1-score (F1), Recall (RR) and Precision (PR) Rates are employed in this work due to class imbalance or high skewness [4] in CASME II and SMIC. With respect to each class " $c$ ", the  $F1_c$ ,  $R_c$  and  $P_c$  are computed from the number of true positive ( $TP_c$ ), false negative ( $FN_c$ ) and false positive ( $FP_c$ ) as follows.

$$RR_c = \frac{TP_c}{TP_c + FN_c} \quad PR_c = \frac{TP_c}{TP_c + FP_c} \quad F1_c = \frac{2R_c P_c}{R_c + P_c}$$

For overall performance on the whole corpus, F1-score, Recall and Precision rates of the dataset are taken as the mean of all respective values of its classes; meanwhile, the overall accuracy (ACC) is computed as  $\frac{\sum_c TP_c}{\sum_c TP_c + FN_c}$ . Table

1 shows average evaluation results across all folds and all classes; meanwhile, Table 3 displays the average evaluation results across all folds for each class. Note that the CASME II paper [3] employs a different evaluation protocol Leave-one-video-out (LOVO); therefore, Yan et al. [3] reports much different and higher accuracy rates for 5-class emotion recognition with the CASME II database in Table 4 of [3]. These results of LOVO are incomparable to those of LOSO in Table 3. As LOVO involves samples belonging to the same subject in both training and testing corpora, it may bias the recognition results.

#### 4.1 Sparse Sampling vs Uniform Sampling

In Table 3, average performances across all LOSO folds and classes of five **SS**, **US**, **US\***, **RA**, and **BL** approaches are compared. Furthermore, Figures 15a, 15c, and 15e show plots of F1-score, Recall Rate and Precision Rate of 5-class subtle emotion recognition on CASME II along percentages ( $p$ ) of preserved frames (45% - 100%) so do Figures 15b, 15d and 15f for the 3-class recognition on SMIC. These plots allow visualization of performances of the approaches with respect to different amounts of redundancy.

Among all evaluation results shown in Figures 15a-15f, the best performances of the recognition system occur at the minimum redundancy  $p = 45\%$  for both CASME II 5-class and SMIC 3-class recognition. Corresponding F1 scores, Precision and Recall rates are listed at the "**SS**" row for both CASME II and SMIC databases. It is experimentally validated that the proposed **SS** approach effectively removes dynamic redundancies compared to any other methods. In accordance with Table 3 and Figures 15a-15f, the **SS** approach improves F1, PR and RR by 40% and ACC by 30% over the second-best performances i.e. **BL** or **US** for CASME II corpus. Meanwhile, F1, RR, PR and ACC of the **SS** approach are improved by 20%-25% over the second best result (of the **US\*** approach) for the SMIC corpus. The performance gap between **SS** and **US** approaches confirms the existence of sparse and redundant dynamics in video samples of micro-expressions. Moreover, the proposed **SS** approach effectively preserves sparse dynamics i.e. emotional expressions, and removes redundant ones i.e. neutral expressions. When redundancies are removed from dynamics of expressions, it increases inter-class and reduces intra-class distances between samples of different emotional classes. As each class consists of video samples from most if not all participating subjects, its intra-class distances are large and inter-class distances are small due to redundant facial identities or neutral expressions. Through elimination of these redundancies, discrimination of video samples, reconstructed by the **SS** approach, w.r.t different emotions is better than that of original samples. Therefore, more compact but meaningful sequences in turn lead to better recognition performances, as shown in Table 3.

The performance of the **US** approach is worse than the **SS** and equivalent to **BL** as shown in Table 3 since TIM arbitrarily removes both discriminative and redundant features by regularly sampling dynamics along a Laplacian manifold. Moreover, the **US** only has better performances than **RA** by a small margin 0.02-0.03 in terms of F1, RR, PR and ACC. The **US** and **US\*** rows show experimental results of

uniform sampling approach with different frame-lengths of generated sequences. Results of these experiments, shown in Figures 15a and 15f, show that the best performance of **US** (green) lines occur at  $p = 60\%$  for CASME II and at  $p = 90\%$  for SMIC. The best F1, PR and RR of the **US** approach are presented at the second row of Table 3. In Yan et al. [3] and Pfister et al. [6], TIM is used for opposite purposes. While CASME II's sequences are lengthened or extrapolated to 150 frames for every sequence in [3], SMIC sequences are shortened or interpolated to 10 frames in [6]. The **US\*** results in Table 3 show recognition performances with respect to those uniform sampling parameters. However, regardless of extrapolation or interpolation, both the **US** and **US\*** approaches are inferior to the proposed **SS** approach in all performance measures, as shown in Table 3 and Figures 15a-15f. Overall, the proposed sparse sampling approach greatly boosts performances of subtle emotion recognition tasks.

#### 4.2 Sparse Sampling vs Other Methods

Yan et al. [3] uses the leave-one-video-out (LOVO) protocol for evaluating performances of 5-class subtle emotion recognition on CASME II database instead of the common LOSO protocol. Though the merit of LOVO evaluation is questionable due to involvement of same subjects in both training and testing phases, we additionally carry out the evaluation based on LOVO for directly comparing the proposed (**SS**) approach with Yan et al. [3] (**YA**). The comparison is valid as all experimental parameters of the proposed approach **SS** for CASME II in Table 5 are also used in [3]. It is noted that the **SS** approach in the table only uses 45% of original frame-lengths. Table 5 shows a significant average improvement in ACC, F1 score, Precision and Recall rates of **SS** over **YA** with the LOVO evaluation protocol. Furthermore, the improvement is also observed in four over five individual classes 'Disgust', 'Happiness', 'Surprise' and 'Repression'. In the 'Others' class, F1 scores of both approaches are equivalent i.e. no improvement of **SS** over **YA**. It is due to skewed sample distribution towards the 'Others' class in CASME II database, shown in Table 2 of [4]. Hence, hyperplanes of SVM classifiers are over-fitted with respect to a cluster of the 'Others' samples. As the biased classifier tends to detect large numbers of true positives as well as false positive samples of the 'Others' class, its recall rate in **YA** is high but precision rate is low. Meanwhile in the **SS** approach, the recall rate drops by 0.03 and precision rate increases by 0.04 from corresponding values of **YA**. **SS** partially removes the bias toward the 'Others' class and improves F1 scores, R and P rates of the other categories.

Besides bench-marking against the baseline method of the CASME II database [3], we also compare the proposed approach to other more recent methods from Huang et al. [10], Oh et al. [11], Liong et al. [12], Wang et al. [33] and Le Ngo et al. [4] with respect to the above evaluation protocols. Table 6 compares performances of the proposed "Sparse Sampling" method with other mentioned state-of-the-art methods. Methods of [11], [12] and [4] are re-implemented; hence, their results in Table 6 are re-evaluated according to parameters of the proposed method with the LOSO protocol. Meanwhile, Huang et al. [10] provides confusion tables for both recognition evaluation on CASME II and



TABLE 5: Leave-one-video-out evaluation on CASME II database of Sparse Sampling (SS) approach and Yan et al. [3] (YA)

	Average				Others (O)			Disgust (D)		
	ACC	F1	RR	PR	F1	RR	PR	F1	RR	PR
SS	.72	.71	.70	.72	.75	.76	.74	.74	.72	.75
YA	.64	.59	.58	.60	.75	.79	.70	.60	.55	.66
(SS-YA)/SS	+12%	+20%	+21%	+20%	0%	-4%	+4%	+23%	+30%	+14%
	Happiness (H)			Surprise (S)			Repression (R)			
	F1	RR	PR	F1	RR	PR	F1	RR	PR	
SS	.56	.61	.51	.82	.80	.83	.67	.59	.76	
YA	.47	.45	.48	.62	.64	.59	.51	.48	.54	
(SS-YA)/SS	+19%	+33%	+6%	+33%	+25%	+41%	+31%	+23%	+41%	

SMIC; therefore, F1-score, recall and precision rates are computed directly from the tables without re-implementation. In general, the proposed approach outperforms most of state-of-the-art methods for both CASME II and SMIC databases except Huang et al. [10] on the CASME II database. Furthermore, "Sparse Sampling" is the only available method that aims to select significant temporal dynamics before feature extraction; meanwhile, the other methods solely focus on designing better features and classifiers. Experimental results shows that proper selection of facial dynamics significantly improves recognition rates of subtle expressions while reducing computational efforts.

TABLE 6: Performance comparison in 5-class recognition for CASME II corpus and 3-class recognition for SMIC against other state-of-the-art methods with the LOSO protocol

	CASME II				SMIC			
	ACC	F1	R	P	ACC	F1	R	P
<b>Sparse Sampling</b>	<b>.49</b>	<b>.51</b>	<b>.47</b>	<b>.55</b>	<b>.58</b>	<b>.60</b>	<b>.60</b>	<b>.60</b>
Huang et al. [10]	.59	.57	.51	.65	.57	.58	.58	.59
Oh et al. [11]	.46	.43	.35	.55	.34	.35	.35	.34
Liong et al. [12]	.42	.38	.36	.41	.53	.54	.55	.53
Wang et al. [33]	.46	.38	.32	.47	.38	.39	.40	.38
Le et al. [4]	.44	.33	.53	.29	.44	.47	.74	.40
Yan et al. [3]	.38	.35	.34	.36	N/A	N/A	N/A	N/A
Pfister et al. [6]	N/A	N/A	N/A	N/A	.40	.40	.40	.41

## 5 CONCLUSION

This work is the first ever endeavour to analyze the dynamics of spontaneous subtle emotions and to learn their sparse structures. Knowledge of these subtle dynamics and the sparsity of data enable spontaneous micro-expressions to be described more discriminately through pre-processing video samples. Both sparse (DMDSP) and uniform sampling (TIM) are compared in theory as principles for removing dynamic redundancy and then both are analyzed in the temporal and spectral domains. The experimental analyses show that the sparse sampling approach is more accurate and consistent than the uniform counterpart. Moreover, when compared against the other state-of-the-art methods in recognizing spontaneous subtle emotion, the performances of the proposed method are very competitive against the state-of-the-art e.g. the best for the SMIC database and the second for CASME II. Therefore, these experimental results confirm the existence of redundancies in the dynamics of spontaneous subtle emotions, and their removals by sparse sampling cause micro-expressions to be more distinctive and recognizable.

## 6 ACKNOWLEDGEMENTS

The authors (ACLN, JS, RP) thank Sébastien Marcel for suggesting to use the F1 score instead of accuracy as a measure, Norman Poh for sharing about his biometrics work involving DMD, and Jonathon Chambers for discussions on DMD and DMDSP. The research and collaboration discussions with Sébastien, Norman and Jonathon were funded by TM (Telekom Malaysia) under the projects UbeAware (MMUE/130152) and 2beAware (MMUE/140098). RP gratefully acknowledges the support by the UK Engineering & Physical Sciences Research Council (EPSRC) under the project Signal Processing Solutions for the Networked Battlespace (EP/K014307/1~2) and the MoD University Defence Research Collaboration (UDRC) in Signal Processing.

## REFERENCES

- [1] R. W. Picard, "Affective computing: from laughter to ieeee," *Affective Computing, IEEE Transactions on*, vol. 1, no. 1, pp. 11–17, 2010.
- [2] C. Shan, S. Gong, and P. W. McOwan, "Facial expression recognition based on local binary patterns: A comprehensive study," *Image and Vision Computing*, vol. 27, no. 6, pp. 803–816, 2009.
- [3] W.-J. Yan, X. Li, S.-J. Wang, G. Zhao, Y.-J. Liu, Y.-H. Chen, and X. Fu, "Casme ii: An improved spontaneous micro-expression database and the baseline evaluation," *PLoS One*, vol. 9, no. 1, p. e86041, 2014.
- [4] A. C. Le Ngo, R. C.-W. Phan, and J. See, "Spontaneous subtle expression recognition: Imbalanced databases & solutions," in *Asian Conference on Computer Vision (ACCV)*, 2014.
- [5] Paul Ekman Group LLC, "Micro expressions," <http://www.paulekman.com/micro-expressions/>, 2016.
- [6] T. Pfister, X. Li, G. Zhao, and M. Pietikainen, "Recognising spontaneous facial micro-expressions," in *Computer Vision (ICCV), 2011 IEEE International Conference on*, 2011, pp. 1449–1456.
- [7] P. J. Schmid, "Dynamic mode decomposition of numerical and experimental data," *Journal of Fluid Mechanics*, vol. 656, pp. 5–28, 2010.
- [8] M. R. Jovanović, P. J. Schmid, and J. W. Nichols, "Sparsity-promoting dynamic mode decomposition," *Physics of Fluids (1994-present)*, vol. 26, no. 2, p. 024103, 2014.
- [9] X. Li, T. Pfister, X. Huang, G. Zhao, and M. Pietikainen, "A spontaneous micro-expression database: Inducement, collection and baseline," in *Automatic Face and Gesture Recognition (FG), 10th IEEE Int. Conf. and Workshops on*, 2013, pp. 1–6.
- [10] X. Huang, S.-J. Wang, G. Zhao, and M. Pietikainen, "Facial micro-expression recognition using spatiotemporal local binary pattern with integral projection," *ICCV Workshop on Computer Vision for Affective Computing, (CV4AC)*, 2015, 2015.
- [11] Y.-H. Oh, A. C. Le Ngo, J. See, S.-T. Liong, R. C.-W. Phan, and H.-C. Ling, "Monogenic riesz wavelet representation for micro-expression recognition," in *Digital Signal Processing, IEEE Int. Conf. on*, 2015, pp. 1237–1241.
- [12] S.-T. Liong, J. See, R. C.-W. Phan, A. C. Le Ngo, Y.-H. Oh, and K. Wong, "Subtle expression recognition using optical strain weighted features," in *1st Workshop on Computer Vision for Affective Computing, Asian Conf. on Computer Vision (ACCV)*, 2014, 2014.
- [13] M. Shreave, S. Godavarthy, D. Goldgof, and S. Sarkar, "Macro-and micro-expression spotting in long videos using spatio-temporal strain," in *Automatic Face & Gesture Recognition and Workshops (FG 2011), 2011 IEEE International Conference on*. IEEE, 2011, pp. 51–56.
- [14] Y. O. S. Polikovsky, Y. Kameda, "Facial micro-expressions recognition using high speed camera and 3d-gradient descriptor," *IET Conference Proceedings*, pp. 16–16(1), January 2009.
- [15] G. Warren, E. Schertler, and P. Bull, "Detecting deception from emotional and unemotional cues," *Journal of Nonverbal Behavior*, vol. 33, no. 1, pp. 59–69, 2009.
- [16] S.-J. Wang, W.-J. Yan, G. Zhao, X. Fu, and C.-G. Zhou, "Micro-expression recognition using robust principal component analysis and local spatiotemporal directional features," in *ECCV 2014 Workshop Spontaneous Facial Behavior Analysis (SFBA)*, 2014, pp. 325–338.
- [17] A.-C. Le-Ngo, Y.-H. Oh, R. C.-W. Phan, and J. See, "Eulerian emotion magnification for subtle expression recognition," in *Acoustics, Speech and Signal Processing (ICASSP), 2016 IEEE Int. Conf. on*, 2016.



- [18] H.-Y. Wu, M. Rubinstein, E. Shih, J. V. Gutttag, F. Durand, and W. T. Freeman, "Eulerian video magnification for revealing subtle changes in the world." *ACM Trans. Graph.*, vol. 31, no. 4, p. 65, 2012.
- [19] N. Wadhwa, M. Rubinstein, F. Durand, and W. T. Freeman, "Phase-based video motion processing," *ACM Transactions on Graphics (TOG)*, vol. 32, no. 4, p. 80, 2013.
- [20] C. Liu, A. Torralba, W. T. Freeman, F. Durand, and E. H. Adelson, "Motion magnification," *ACM Transactions on Graphics (TOG)*, vol. 24, no. 3, pp. 519–526, 2005.
- [21] J. Grosek and J. N. Kutz, "Dynamic mode decomposition for real-time background/foreground separation in video," *arXiv preprint arXiv:1404.7592*, 2014.
- [22] G. Zhao and M. Pietikainen, "Dynamic texture recognition using local binary patterns with an application to facial expressions," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 29, no. 6, pp. 915–928, 2007.
- [23] M. Shreve, S. Godavathy, V. Manohar, D. Goldgof, and S. Sarkar, "Towards macro-and micro-expression spotting in video using strain patterns," in *Applications of Computer Vision (WACV), 2009 Workshop on*. IEEE, 2009, pp. 1–6.
- [24] Y.-J. Liu, J.-K. Zhang, W.-J. Yan, S.-J. Wang, G. Zhao, and X. Fu, "A main directional mean optical flow feature for spontaneous micro-expression recognition," *Affective Computing, IEEE Transactions on*, vol. PP, no. 99, pp. 1–1, 2015.
- [25] X. Huang, Q. He, X. Hong, G. Zhao, and M. Pietikainen, "Improved spatiotemporal local monogenic binary pattern for emotion recognition in the wild," in *Proc. of the 16th Int. Conf. on Multimodal Interaction*, 2014, pp. 514–520.
- [26] Y.-H. Oh, A.-C. Le-Ngo, R. C.-W. Phan, J. See, and H.-C. Ling, "Intrinsic two-dimensional local structures for micro-expression recognition," in *Acoustics, Speech and Signal Processing (ICASSP), 2016 IEEE Int. Conf. on*, 2016.
- [27] Z. Zhou, G. Zhao, and M. Pietikainen, "Towards a practical lipreading system," in *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conf. on*. IEEE, 2011, pp. 137–144.
- [28] S. Yan, D. Xu, B. Zhang, H.-J. Zhang, Q. Yang, and S. Lin, "Graph embedding and extensions: a general framework for dimensionality reduction," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 29, no. 1, pp. 40–51, 2007.
- [29] X. He, D. Cai, S. Yan, and H.-J. Zhang, "Neighborhood preserving embedding," in *Computer Vision, 2005. ICCV 2005. Tenth IEEE International Conference on*, vol. 2. IEEE, 2005, pp. 1208–1213.
- [30] Z. Wen, D. Goldfarb, and W. Yin, "Alternating direction augmented lagrangian methods for semidefinite programming," *Mathematical Programming Computation*, vol. 2, no. 3-4, pp. 203–230, 2010.
- [31] P. Lucey, J. F. Cohn, T. Kanade, J. Saragih, Z. Ambadar, and I. Matthews, "The Extended Cohn-Kanade Dataset (CK+): A complete dataset for action unit and emotion-specified expression," in *Computer Vision and Pattern Recognition Workshops (CVPRW), IEEE Int. Conf. on*, 2010, pp. 94–101.
- [32] G. Littlewort, M. S. Bartlett, I. Fasel, J. Susskind, and J. Movellan, "Dynamics of facial expression extracted automatically from video," *Image and Vision Computing*, vol. 24, no. 6, pp. 615–625, 2006.
- [33] Y. Wang, J. See, R. C.-W. Phan, and Y.-H. Oh, "Efficient spatiotemporal local binary patterns for spontaneous facial micro-expression recognition." *PloS one*, vol. 10, no. 5, p. e0124674, 2015.



**Anh-Cat Le-Ngo, Ph.D.** was born in Ho Chi Minh city, Vietnam. He has finished his doctorate training in Computer Vision and Image Processing in University of Nottingham in 2015. Since 2013, he has been working as a post-doc in the UbeAware and 2beAware projects under the guidance of Prof. Raphael C.-W. Phan. His main research interests are computer vision systems (i.e. micro-expression recognition system (MERS), advance driving assistance system (ADAS), etc), signal & image processing mathematics (i.e. monogenic signal, analytic signal, Hilber & Riesz transforms, Dynamic Mode Decomposition, etc) and machine learning (i.e. Long Short Term Memory (LSTM), Recurrent Neural Network (RNN), Selective Transfer Machine (STM), etc).



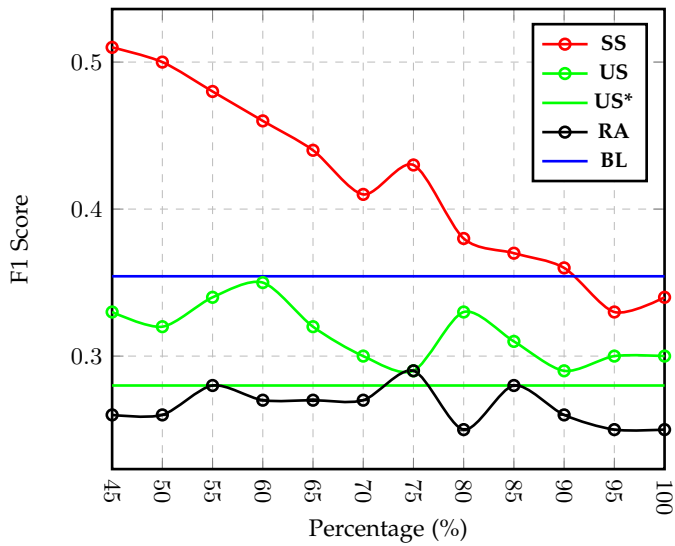
**John See, Ph.D.** received his PhD in Computer Science, MEngSc and BEng degrees from Multimedia University (MMU), Malaysia. He is currently working as a Senior Lecturer and leader of the Pattern Recognition sub-cluster of the Centre for Visual Computing at Multimedia University, Malaysia. His research interests covers a diverse range of topics in computer vision and pattern recognition, particularly in the domain of video-based biometrics, visual surveillance, affective computing, image aesthetics and deep

learning. His passion remains in designing effective and efficient algorithms for visual recognition tasks.

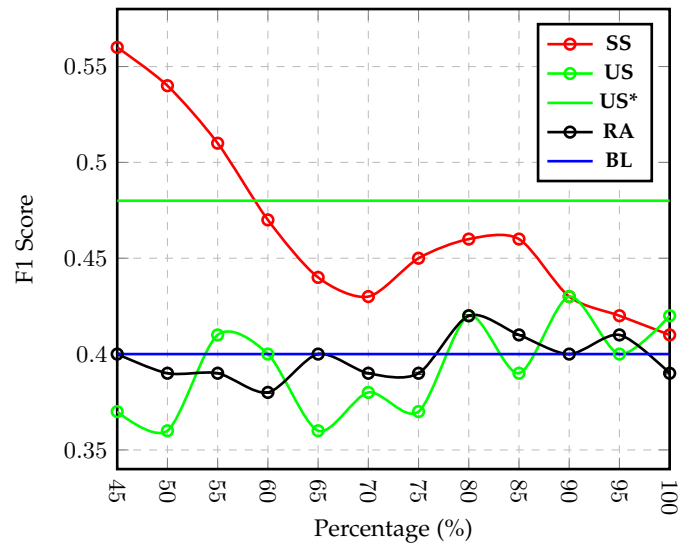


**Raphael C.-W. Phan, Ph.D.** holds the full chair in security engineering at the Faculty of Engineering, Multimedia University (MMU). Prior to joining MMU, he worked at British, Swiss and Australian universities. He has led projects funded by the UK Engineering & Physical Sciences Research Council (EPSRC), UK Ministry of Defence (MoD), as well as the Malaysian government and industry. Raphael's research passion spans the breadth of security & privacy, from cryptography and security protocols

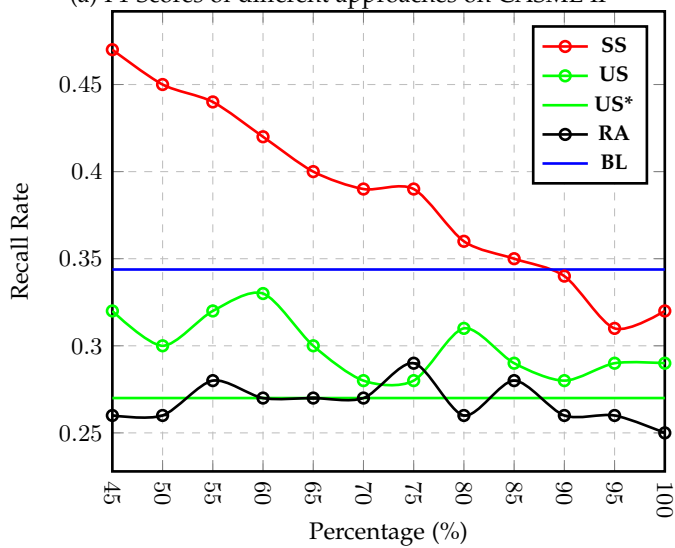
to subtle and/or hidden emotions: essentially aiming to see beyond the unseen. He was co-designer of the BLAKE hash function, one of the five finalists of the US National Institute of Standards & Technology (NIST) SHA-3 Hash Function Competition. Raphael is regularly invited to serve in the technical program committees of peer-reviewed security conferences, and is Program Chair of Mycrypt 2016, the first such conference of its kind with focus on malicious and/or out-of-the-box cryptographic paradigms.



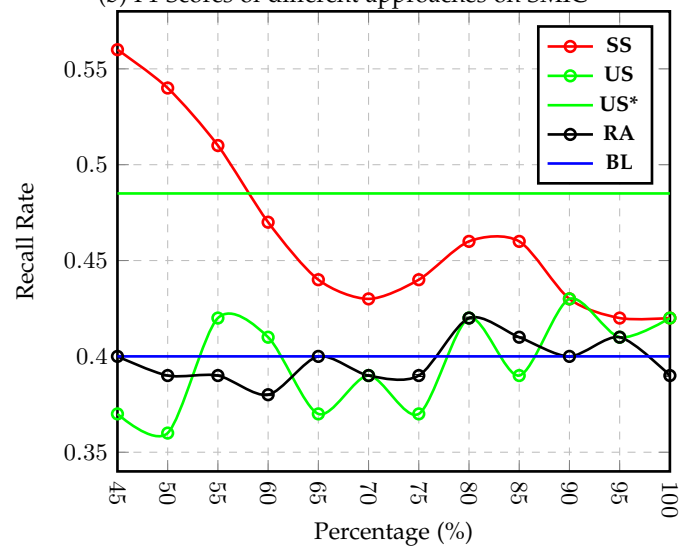
(a) F1 Scores of different approaches on CASME II



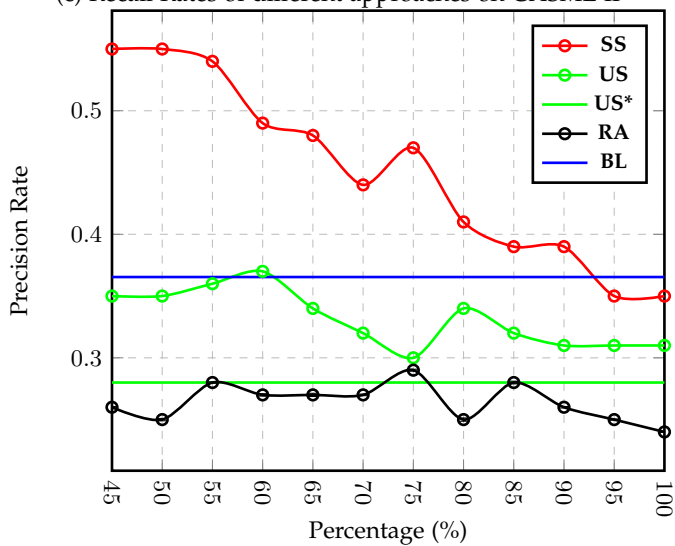
(b) F1 Scores of different approaches on SMIC



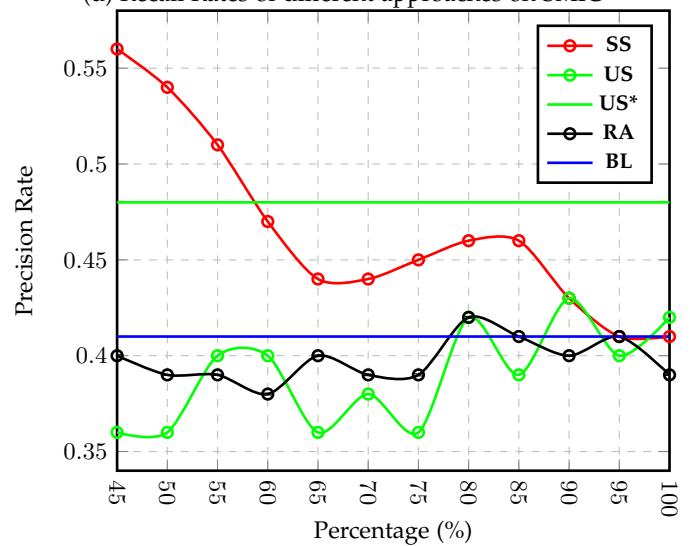
(c) Recall Rates of different approaches on CASME II



(d) Recall Rates of different approaches on SMIC



(e) Precision Rates of different approaches on CASME II



(f) Precision Rates of different approaches on SMIC

Fig. 15: Performance metrics (F1 score, Recall and Precision Rates) of various temporally sampling methods with respect to a range of percentages of preserved frames (45%-100% original frame-length) on CASME II and SMIC databases