

A Branch-and-Bound Framework for Unsupervised Common Event Discovery

Wen-Sheng Chu¹  · Fernando De la Torre¹ ·
Jeffrey F. Cohn^{1,2} · Daniel S. Messinger³

Received: 3 June 2016 / Accepted: 12 January 2017 / Published online: 9 February 2017
© Springer Science+Business Media New York 2017

Abstract Event discovery aims to discover a temporal segment of interest, such as human behavior, actions or activities. Most approaches to event discovery within or between time series use supervised learning. This becomes problematic when relevant event labels are unknown, are difficult to detect, or not all possible combinations of events have been anticipated. To overcome these problems, this paper explores Common Event Discovery (CED), a new problem that aims to discover common events of variable-length segments in an unsupervised manner. A potential solution to CED is searching over all possible pairs of segments, which would incur a prohibitive quartic cost. In this paper, we propose an efficient branch-and-bound (B&B) framework that avoids exhaustive search while guaranteeing a globally optimal solution. To this end, we derive novel bounding functions for various

commonality measures and provide extensions to multiple commonality discovery and accelerated search. The B&B framework takes as input any multidimensional signal that can be quantified into histograms. A generalization of the framework can be readily applied to discover events at the same or different times (synchrony and event commonality, respectively). We consider extensions to video search and supervised event detection. The effectiveness of the B&B framework is evaluated in motion capture of deliberate behavior and in video of spontaneous facial behavior in diverse interpersonal contexts: interviews, small groups of young adults, and parent-infant face-to-face interaction.

Keywords Common event discovery · Synchrony discovery · Video indexing · Event detection · Human interaction · Unsupervised learning · Global optimization · Branch and bound · Bag-of-words

Communicated by Shin'ichi Satoh.

Electronic supplementary material The online version of this article (doi:[10.1007/s11263-017-0989-7](https://doi.org/10.1007/s11263-017-0989-7)) contains supplementary material, which is available to authorized users.

✉ Wen-Sheng Chu
wschu@cmu.edu

Fernando De la Torre
ftorre@cs.cmu.edu

Jeffrey F. Cohn
jeffcjohn@pitt.edu

Daniel S. Messinger
dmessinger@miami.edu

¹ Robotics Institute, Carnegie Mellon University, Pittsburgh, PA, USA

² Department of Psychology, University of Pittsburgh, Pittsburgh, PA, USA

³ Department of Psychology, University of Miami, Coral Gables, FL, USA

1 Introduction

Event detection is a central topic in computer vision. Most approaches to event detection use one or another form of supervised learning. Labeled video from experts or naive annotators is used as training data, classifiers are trained, and then used to detect individual occurrences or pre-defined combinations of occurrences in new video. While supervised learning has well-known advantages for event detection, limitations might be noted. One, because accuracy scales with increases in the number of subjects for whom annotated video is available, sufficient numbers of training subjects are essential (Girard et al. 2015; Chu et al. 2016). With too few training subjects, supervised learning is under-powered. Two, unless an annotation scheme is comprehensive, important events may go unlabeled, unlearned, and ultimately

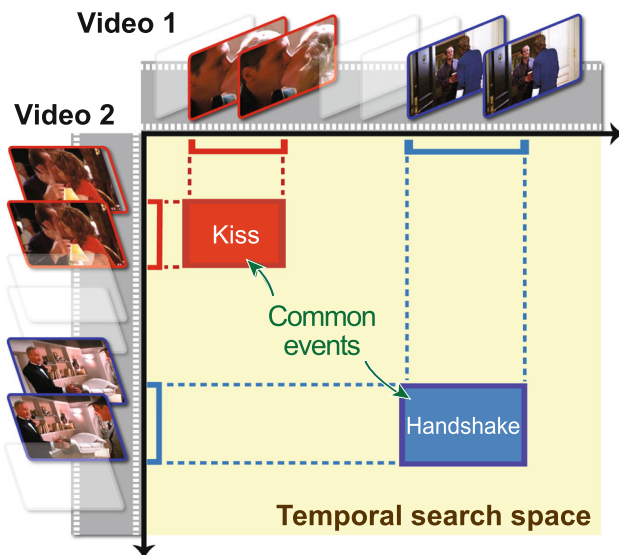


Fig. 1 An illustration of Common event discovery (CED). Given two videos, common events (*kiss* and *handshake*) of different lengths in the two videos are discovered in an unsupervised manner

undetected. Three and perhaps most important, discovery of similar or matching events is limited to combinations of actions that have been specified in advance. Unanticipated events go unnoticed. To enable the discovery of novel recurring or matching events or patterns, unsupervised discovery is a promising option.

To detect recurring combinations of actions without labels, this paper addresses Common event discovery (CED), a relatively unexplored problem that discovers common temporal events in variable-length segments in an unsupervised manner. The goal of CED is to detect pairs of segments that retain maximum visual commonality. CED is fully unsupervised, so no prior knowledge about events is required. We need not know what the common events are, how many there are, or when they may begin and end. Figure 1 illustrates the concept of CED for video. In an exhaustive search of variable-length video segments, *kissing* and *handshake* event matches are discovered between videos.

A naive approach to CED would be to use a sliding window. That is, to exhaustively search all possible pairs of temporal segments and select pairs that have the highest similarities. Because the complexity of sliding window methods is quartic with the length of video, *i.e.*, $\mathcal{O}(m^2n^2)$ for two videos of lengths m and n , this cost would be computationally prohibitive in practice. Even in relatively short videos of 200 and 300 frames, there would be in excess of *three billion* possible matches to evaluate at different lengths and locations.

To meet the computational challenge, we propose to extend the Branch-and-Bound (B&B) method for CED. For supervised learning, B&B has proven an efficient technique to detect image patches (Lampert et al. 2009) and video volumes (Yuan et al. 2011). Because previous bounding

functions of B&B are designed for supervised detection or classification, which require pre-trained models, previous B&B methods could not be directly applied to CED. For this reason, we derive novel bounding functions for various commonality measures, including ℓ_1/ℓ_2 distance, intersection kernel, χ^2 distance, cosine similarity, symmeterized cross entropy, and symmeterized KL-divergence.

For evaluation, we apply the proposed B&B to application of discovering events at the same or different times (synchrony and event commonality, respectively), and variable-length segment-based event detection. We conduct the experiments on three datasets of increasing complexity: Posed motion capture and unposed, spontaneous video of mothers and their infants and of young adults in small groups. We report distance and similarity metrics and compare discovery with expert annotations. Our main contributions are:

1. *A new CED problem* Common event discovery (CED) in video is a relatively unexplored problem in computer vision. Results indicate that CED achieves moderate convergence with supervised approaches, and is able to identify novel patterns both within and between time series.
2. *A novel, unsupervised B&B framework* With its novel bounding functions, the proposed B&B framework is computationally efficient and entirely general. It takes any signals that can be quantified into histograms and with minor modifications adapts readily to diverse applications. We consider four: common event discovery, synchronous event discovery, video search, and supervised segment-based event detection.

A preliminary version of this work appeared as Chu et al. (2012, 2015). In this paper, we integrate these two approaches with video search and supervised segment-based event detection, and provide a principal way of deriving bounding functions in the new, unsupervised framework. We also present new experiments on supervised event detection with comparisons to alternative methods. The rest of this paper is organized as follows. Section 2 discusses related work. Section 3 presents the proposed B&B framework for common event discovery. Section 4 applies the framework to tasks of varying complexity. Section 5 extends the B&B framework to discovery among more than two videos and considers acceleration using warm-start strategy and parallelism. Section 6 provides evaluation on unsupervised and supervised tasks with unsynchronous and synchronous videos. Section 7 concludes the paper with future work.

2 Related Work

This paper is closely related to event detection methods, and unsupervised discovery in images and videos. Below we review each in turn.

2.1 Event Detection

CED closely relates to event detection. Below we categorize prior art into supervised and unsupervised approaches, and discuss each in turn.

Supervised event detection Supervised event detection is well-developed in computer vision. Events can be defined as temporal segments that involve either a single pattern of interest or an interaction between multiple patterns. For single-pattern event detection, popular examples include facial expression recognition (Sangineto et al. 2014; Littlewort et al. 2006; Valstar and Pantic 2006; Lucey et al. 2010; Du et al. 2014), surveillance system (Feris et al. 2014), activity recognition (Gao et al. 2015; Yang et al. 2013a, b; Reddy and Shah 2013; Duchenne et al. 2009; Jhuang et al. 2007), and sign language recognition (Cooper and Bowden 2009). These approaches aim to detect a temporal pattern that associates with a pre-defined human behavior, action, or activity.

Events may also be defined as the co-occurrence of discrete actions or activities. For instance, (Brand et al. 1997) treated each arm as a process, and proposed to recognize gestures by modeling motion trajectories between multiple processes using coupled hidden Markov models (CHMMs). Following up, (Oliver et al. (2000)) proposed a CHMM-based system, with pedestrian trajectories, to detect and recognize interactions between people, such as following another person, altering one's path to encounter another, etc. (Hongeng and Nevatia (2001) proposed a hierarchical trajectory representation along with a temporal logic network to address complex interactions such as a "stealing" scenario. More recently, (Liu et al. (2010) proposed to recognize group behavior in AAL environment (nursing homes), considering a switch control module that alternates between two HMM-based methods built on motion and poses of individuals. (Messinger et al. (2010) focused on specific annotated social signals, i.e., smiling and gaze, and characterized the transition between behavior states by a maximum likelihood approach. Interested readers are referred to (Chaaraoui et al. (2012) for a review. These techniques, however, require adequate labeled training data, which can be time-consuming to collect and not always available.

Unsupervised event detection The closest to our study is unsupervised approaches that require no annotations. For instance, (Zheng et al. (2011) presented a coordinated motion model to detect motion synchrony in a group of individuals such as fish schools and bird flocks. (Zhou et al. (2013) proposed aligned cluster analysis that extended spectral clustering to cluster time series, and applied the technique to discover facial events in unsupervised manner. On the other hand, *time series motifs*, defined as the closest pair of subsequences in one time series stream, can be discovered with a tractable exact algorithm (Mueen and Keogh 2010), or an approximated algorithm that is capable of tackling never-

ending streams (Begum and Keogh 2014). Some attempts at measuring interactional synchrony include using face tracking and expressions (Yu et al. 2013), and rater-coding and pixel changes between adjacent frames (Schmidt et al. 2012). (Nayak et al. (2012) presented iterated conditional modes to find most recurrent sign in all occurrences of sign language sentences.

Common events refer to two or more actions that are similar either in form or in timing. The meaning of similarity depends upon the choice of features, similarity metrics, and the threshold to accept similarity. While cluster analysis or mode finding could be considered a potential method, it is not well-suited for common event discovery for some reasons. First, cluster analysis and mode finding methods are designed for discovering the instances or values that appear most often; yet, common events could appear rarely. Second, cluster analysis and mode finding methods consider all instances to obtain statistical "groups" or "modes"; common events are a sparse subset of instances with high similarity. Finally, cluster analysis and mode finding methods for time series require temporal segmentation as a pre-processing procedure; common event discovery has no such requirement.

2.2 Unsupervised Discovery

For static images, unsupervised discovery of re-occurring patterns has proven informative, driven by wide applications in co-segmentation (Chu et al. 2010; Liu and Yan 2010; Mukherjee et al. 2011), grammar learning (Zhu and Mumford 2006), irregularity detection (Boiman and Irani 2005) and automatic tagging (Schindler et al. 2008) have been driving forces. Discovery of common patterns in videos is a relatively unexplored problem. See (Wang et al. (2014) for a survey.

For video, to our best knowledge, this study is the first to discover common events in an unsupervised manner. Our work is inspired by recent success on using B&B for efficient search. (Lampert et al. (2009) proposed Efficient Subwindow Search (ESS) to find the optimal subimage that maximizes the Support Vector Machine score of a pre-trained classifier. (Hoai et al. (2011) combine SVM with dynamic programming for efficient temporal segmentation. (Yuan et al. (2011) generalized Lampert's 4-D search to the 6-D Spatio-Temporal Branch-and-Bound (STBB) search by incorporating time, to search for spatio-temporal volumes. However, unlike CED, these approaches are supervised and require a training stage.

Recently, there have been interests on temporal clustering algorithms for unsupervised discovery of human actions. (Wang et al. (2006) used deformable template matching of shape and context in static images to discover action classes. (Si et al. (2011) learned an event grammar by clustering event co-occurrence into a dictionary of atomic actions. (Zhou et al. (2010) combined spectral clustering and dynamic time warping to cluster time series, and applied it to learn tax-

onomies of facial expressions. Turaga et al. (2009) used extensions of switching linear dynamical systems for clustering human actions in video sequences. However, if we cluster two sequences that each has only one segment in common, previous clustering methods would likely need many clusters to find the common segments. In our case, CED focuses only on common segments and avoids clustering all video segments, which is computationally expensive and prone to local minimum.

Another unsupervised technique related to CED is motif detection (Minnen et al. 2007; Mueen and Keogh 2010). Time series motif algorithms find repeated patterns within a single sequence. Minnen et al. (2007) discovered motifs as high-density regions in the space of all subsequences. Mueen and Keogh (2010) further improved the motif discovery problem using an online technique, maintaining the exact motifs in real-time performance. Nevertheless, these work detects motifs within only one sequence, but CED considers two (or more) sequences. Moreover, it is unclear how these technique can be robust to noise.

Finally, CED is also related to the longest common subsequence (LCS) (Gusfield 1997; Maier 1978; Paterson and Dančik 1994). The LCS problem consists on finding the longest subsequence that is common to a set of sequences (often just two) (Paterson and Dančik 1994; Wang and Velipasalar 2009). Closer to our work is the algorithm for discovering longest consecutive common subsequence (LCCS) (Wang and Velipasalar 2009), which finds the longest contiguous part of original sequences (e.g., videos). However, different from CED, these approaches have a major limitation in that they find only identical subsequences, and hence are sensitive to noisy signals in realistic videos.

3 A Branch-and-Bound Framework for Common Event Discovery (CED)

This section describes our representation of time series, a formulation of CED, the proposed B&B framework, and the newly derived bounding functions that fit into the B&B framework.

3.1 Representation of Time Series

Bag of Temporal Words (BoTW) model (Sivic and Zisserman 2003; Yuan et al. 2011) has been shown effective in many video analysis problems, such as action recognition (Brendel and Todorovic 2011; Han et al. 2009; Laptev et al. 2008; Liu et al. 2011; Sadanand and Corso 2012). This section modifies the BoTW model to describe the static and dynamic information of a time series. Suppose a time series \mathbf{S} can be described as a set of feature vectors $\{\mathbf{x}_j\}$ for each frame j

(see notation¹). For instance, a feature vector can be facial shape in face videos or joint angles in motion capture videos. Given such features, we extract two types of information: *observation info* from a single frame, and *interaction info* from two consecutive frames. Denote $\mathbf{S}[b, e] = \{\mathbf{x}_j\}_{j=b}^e$ as a temporal segment between the b -th and the e -th frames, we consider a segment-level feature mapping:

$$\varphi_{\mathbf{S}[b,e]} = \sum_{j=b}^e \begin{bmatrix} \phi^{\text{obs}}(\mathbf{x}_j) \\ \phi^{\text{int}}(\mathbf{x}_j) \end{bmatrix}. \quad (1)$$

The observation info $\phi^{\text{obs}}(\mathbf{x}_j)$ describes the “pseudo” probability of \mathbf{x}_j belonging to a latent state, and the interaction info $\phi^{\text{int}}(\mathbf{x}_j)$ describes transition probability of states between two consecutive frames. To obtain $\phi^{\text{obs}}(\mathbf{x}_j)$, we performed k -means to find K centroids $\{\mathbf{c}_k\}_{k=1}^K$ as the hidden states. Then, we computed $\phi^{\text{obs}}(\mathbf{x}_j) \in [0, 1]^K$ with the k -th element computed as $\exp(-\gamma \|\mathbf{x}_j - \mathbf{c}_k\|^2)$ and γ chosen as an inverse of the median distance of all samples to the centroids. An interaction info $\phi^{\text{int}}(\mathbf{x}_j) \in [0, 1]^{K^2}$ is computed as:

$$\phi^{\text{int}}(\mathbf{x}_j) = \text{vec} \left(\phi^{\text{obs}}(\mathbf{x}_j) \otimes \phi^{\text{obs}}(\mathbf{x}_{j+1}) \right), \quad (2)$$

where \otimes denotes a Kronecker product of two observation vectors. As a result, each temporal segment is represented as an ℓ_2 -normalized feature vector of dimension $(K^2 + K)$.

Because this representation accepts almost arbitrary features, any signal, even with negative values, that can be quantified into histograms can be directly applied. One notable benefit of the histogram representation is that it allows for fast recursive computation using the concept of *integral image* (Viola and Jones 2004). That is, the segment-level representation for $\mathbf{S}[b, e]$ can be computed as $\varphi_{\mathbf{S}[b,e]} = \varphi_{\mathbf{S}[1,e]} - \varphi_{\mathbf{S}[1,b-1]}$, which only costs $\mathcal{O}(1)$ per evaluation. Based on the time series representation, we develop our approach below.

3.2 Problem Formulation

To establish notion, we begin with two time series \mathbf{S}^1 and \mathbf{S}^2 with m and n frames respectively. The goal of common event discovery (CED) is to find two temporal segments with intervals $[b_1, e_1] \subseteq [1, m]$ and $[b_2, e_2] \subseteq [1, n]$ such that their visual commonality is maximally preserved. We formulate CED:

$$\boxed{\text{CED}} \quad \max_{\{b_1, e_1, b_2, e_2\}} f \left(\varphi_{\mathbf{S}^1[b_1, e_1]}, \varphi_{\mathbf{S}^2[b_2, e_2]} \right),$$

¹ Bold capital letters denote a matrix \mathbf{X} , bold lower-case letters a column vector \mathbf{x} . \mathbf{x}_i represents the i th column of the matrix \mathbf{X} . x_{ij} denotes the scalar in the i th row and j th column of the matrix \mathbf{X} . All non-bold letters represent scalars.

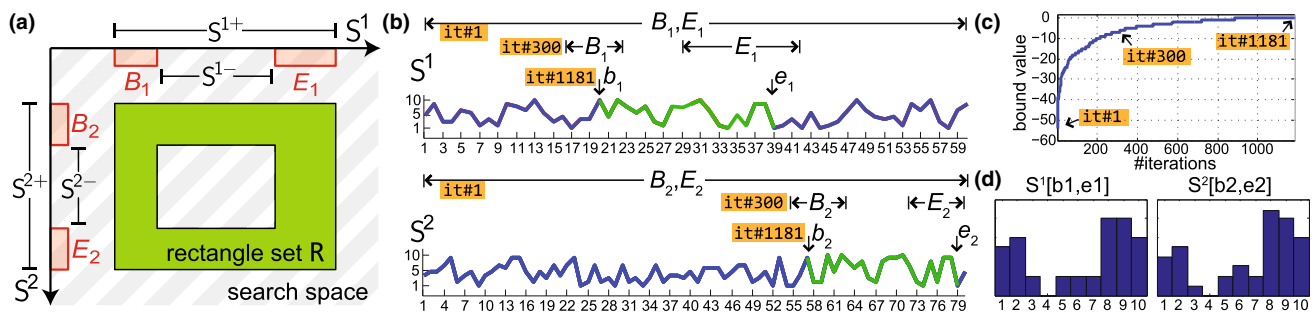


Fig. 2 An example of CED on two 1-D time series: **a** an illustration of our notation (see Sect. 3.3). **b** Searching intervals at iterations (*it*) #1, #300 and #1181 over sequences S^1 and S^2 . Commonalities $S^1[b_1, e_1]$ and $S^2[b_2, e_2]$ are discovered at convergence (#1181). **c** Convergence

curve w.r.t. bounding value and *it*. **d** Histograms of the discovered commonalities. In this example, a naive sliding window approach needs more than 5 million evaluations, while the proposed B&B method converges at iteration 1181 using $\ell = 20$

$$\text{subject to } \ell \leq e_i - b_i, \quad \forall i \in \{1, 2\}, \quad (3)$$

where $f(\cdot, \cdot)$ is a commonality measure between two time series representations, and ℓ controls the minimal length for each temporal segment to avoid a trivial solution. More details about $f(\cdot, \cdot)$ are discussed in Sect. 3.4. Problem (3) is non-convex and non-differentiable, and thus standard convex optimization methods can not be directly applied. A naive solution is an exhaustive search over all possible locations for $\{b_1, e_1, b_2, e_2\}$. However, it leads to an algorithm with computational complexity $\mathcal{O}(m^2n^2)$, which is prohibitive for regular videos with hundreds or thousands of frames. To address this issue, we introduce a branch-and-bound (B&B) framework to efficiently and globally solve (3).

Note that, although ℓ controls the minimal length of discovered temporal segments, the optimal solution can be of length greater than ℓ . For instance, consider two 1-D time series $S^1 = [1, 2, 2, 1]$ and $S^2 = [1, 1, 3]$. Suppose we measure $f(\cdot, \cdot)$ by ℓ_1 distance, where smaller values indicate higher commonality. Let the minimal length $\ell = 3$, and represent their 3-bin histograms as $\varphi_{S^1[1,4]} = [2, 2, 0]$, $\varphi_{S^1[1,3]} = [1, 2, 0]$ and $\varphi_{S^2} = [2, 0, 1]$. Showing the distance $f_{\ell_1}(\varphi_{S^1[1,4]}, \varphi_{S^2}) = 3 < 4 = f_{\ell_1}(\varphi_{S^1[1,3]}, \varphi_{S^2})$, we prove by contradiction.

3.3 Optimization by Branch and Bound (B&B)

With a proper bounding function, B&B has been shown empirically more efficient than straight enumeration. B&B can eliminate regions that provably do not contain an optimal solution. This can be witnessed in many computer vision problems, *e.g.*, object detection (Lampert et al. 2009; Lehmann et al. 2011), video search (Yuan et al. 2011), pose estimation (Sun et al. 2012) and optimal landmark detection (Amberg and Vetter 2011). Inspired by previous success, this section describes the proposed B&B framework that globally solves (3).

Problem interpretation As depicted in Fig. 1, we interpret Problem (3) as searching a rectangle in the 2-D space formed by two time series. A rectangle $r \doteq [b_1, e_1, b_2, e_2]$ in the search space indicates one candidate solution corresponding to $S^1[b_1, e_1]$ and $S^2[b_2, e_2]$. To allow a more efficient representation for searching, we parameterize each step as searching over sets of candidate solutions. That is, we search over *intervals* instead of individual value for each parameter. Each parameter interval corresponds to a rectangle set $R \doteq B_1 \times E_1 \times B_2 \times E_2$ in the search space, where $B_i = [b_i^{lo}, b_i^{hi}]$ and $E_i = [e_i^{lo}, e_i^{hi}]$ ($i \in \{1, 2\}$) indicate tuples of parameters ranging from frame *lo* to frame *hi*. Given the rectangle set *R*, we denote the longest and the shortest possible segments as S^{i+} and S^{i-} respectively. We denote $|R|$ as the number of rectangles in *R*. Figure 2a shows an illustration of the notation.

The B&B framework With the problem interpreted above, we describe here the proposed B&B framework. Algorithm 1 summarizes the procedure. To maintain the search process, we employ a priority queue denoted as *Q*. Each state in *Q* contains a rectangle set *R*, its upper bound $u(R)$ and lower bound $l(R)$. Each iteration starts by selecting a rectangle set *R* from the top state, which is defined as the state containing the minimal upper bound for $f(\cdot, \cdot)$. Given this structure, the algorithm repeats a *branch* step and a *bound* step until *R* contains a unique entry.

In the *branch* step, each rectangle set *R* is split by its largest interval into two disjoint subsets. For example, suppose E_2 is the largest interval, then $R \rightarrow R' \cup R''$ where $E'_2 = [e_2^{lo}, \lfloor \frac{e_2^{lo} + e_2^{hi}}{2} \rfloor]$ and $E''_2 = [\lfloor \frac{e_2^{lo} + e_2^{hi}}{2} \rfloor + 1, e_2^{hi}]$. In the *bound* step, we calculate the bounds for each rectangle set, and then update new rectangle sets and their bounds into *Q*. The computed bounds tell the worst possible values in $f(\cdot, \cdot)$, and therefore enable the algorithm to efficiently discard unlikely rectangle sets where their bounds are worse than the current best. The algorithm terminates when *R* contains a unique entry.

Algorithm 1: Common Event Discovery (CED)

Input : Collection of frame-based features for sequences S^1, S^2 ; minimal length ℓ
Output: The optimal rectangle r^* in the search space

```

1  $Q \leftarrow$  empty priority queue;           // Initialize  $Q$ 
2  $R \leftarrow [1, m] \times [1, m] \times [1, n] \times [1, n]$ ; // Initialize  $R$ 
3  $r^* \leftarrow \text{BnB}(Q, R)$ ; // Obtain the optimal  $r$  using BnB
4 return  $r^*$ ;
1 Procedure  $\text{BnB}(Q, R)$ 
2   while  $|R| \neq 1$  do
3      $R \rightarrow R' \cup R''$ ;           // Branch step
4      $Q.\text{push}(\text{bound}(R'), R')$ ; // Push  $R_1$  and bound
5      $Q.\text{push}(\text{bound}(R''), R'')$ ; // Push  $R_2$  and bound
6      $R \leftarrow Q.\text{pop}()$ ;       // Pop from  $Q$ 
7   end
8   return  $R$ ;
```

a unique entry, i.e., $|R| = 1$. Figure 2b–d show an example of CED for discovering commonality between two 1-D time series. Despite that in the worst case the complexity of B&B can be still $\mathcal{O}(m^2n^2)$, we will experimentally show that in general B&B is much more efficient than naive approaches.

3.4 Construction of Bounding Functions

One crucial aspect of the proposed B&B framework is the novel bounding functions for measuring commonality between two time series. The commonality measures can interchangeably be formed in terms of distance or similarity functions. Below we describe the conditions of bounding functions, and then construct the bounds.

Conditions of bounding functions Recall that R represents a rectangle set and $r \doteq [b_i, e_i, b_j, e_j]$ represents a rectangle corresponding to two subsequences $S^i[b_i, e_i]$ and $S^j[b_j, e_j]$. Without loss of generality, we denote $f(r) = f(\varphi_{S^i[b_i, e_i]}, \varphi_{S^j[b_j, e_j]})$ as the commonality measure between $S^i[b_i, e_i]$ and $S^j[b_j, e_j]$. To harness the B&B framework, we need to find an upper bound $u(R)$ and a lower bound $l(R)$ that bounds the values of f over a set of rectangles. A proper bounding function has to satisfy the conditions:

- (a) $u(R) \geq \max_{r \in R} f(r)$, Bounding conditions
 (b) $l(R) \leq \min_{r \in R} f(r)$,
 (c) $u(R) = f(r) = l(R)$, if r is the only element in R .

Conditions (a) and (b) ensure that $u(R)$ and $l(R)$ appropriately bound all candidate solutions in R from above and from below, whereas (c) guarantees the algorithm to converge to the optimal solution. With both lower and upper bounds, one can further prune the priority queue for speeding the search, i.e., eliminate rectangle sets R' that satisfy $l(R') > u(R)$ (Balakrishnan et al. 1991).

Bound histogram bins Let S^i denote the i -th time series and can be represented as an unnormalized histogram \mathbf{h}^i or a normalized histogram $\hat{\mathbf{h}}^i$ using the representation in Sect. 3.1. Denote h_k^i and \hat{h}_k^i as the k -th bin of \mathbf{h}^i and $\hat{\mathbf{h}}^i$, respectively. The normalized histogram is defined as $\hat{h}_k^i = h_k^i / |S^i|$, where $|S^i| = \sum_k h_k^i$. $\|\mathbf{S}^i\| = \sqrt{\sum_k (h_k^i)^2}$ is the Euclidean norm of histogram of S^i . Considering histograms of S^{i+} and S^{i-} , we can bound their k -th histogram bin:

$$0 \leq h_k^{i-} \leq h_k^i \leq h_k^{i+}, \quad \forall i. \quad (4)$$

Given a rectangle $r = [b_1, e_1, b_2, e_2]$ and denote $\underline{h}_k^i = \frac{h_k^{i-}}{|S^{i+}|}$ and $\overline{h}_k^i = \frac{h_k^{i+}}{|S^{i-}|}$. For normalized histograms, we use the fact that $|S^{i-}| \leq |S^i[b_i, e_i]| \leq |S^{i+}|$. Then we can rewrite (4) for bounding the normalized bins:

$$0 \leq \underline{h}_k^i \leq \hat{h}_k^i \leq \overline{h}_k^i, \quad \forall i. \quad (5)$$

Below we use Eq. (5) to construct bounds for various commonality measures with normalized histograms, whereas those with unnormalized histograms can be likewise obtained.

Bound commonality measures: Given two time series S^i and S^j represented as normalized histograms $\hat{\mathbf{h}}^i$ and $\hat{\mathbf{h}}^j$ respectively, we provide bounding functions for various commonality measures: ℓ_1/ℓ_2 distance, histogram intersection, χ^2 distance, cosine similarity, symmetrized KL divergence, and symmetrized cross entropy. These measures have been widely applied to many tasks such as object recognition (Everingham et al. 2006; Lampert et al. 2009) and action recognition (Brendel and Todorovic 2011; Han et al. 2009; Laptev et al. 2008; Liu et al. 2011; Sadanand and Corso 2012).

1. ℓ_1/ℓ_2 distance Applying the min/max operators on (4), we get

$$\begin{aligned} \min(h_k^{i-}, h_k^{j-}) &\leq \min(h_k^i, h_k^j) \leq \min(h_k^{i+}, h_k^{j+}), \\ \text{and } \max(h_k^{i-}, h_k^{j-}) &\leq \max(h_k^i, h_k^j) \leq \max(h_k^{i+}, h_k^{j+}). \end{aligned} \quad (6)$$

Reordering the inequalities, we obtain the upper bound u_k and lower bound l_k for the k -th histogram bin:

$$\begin{aligned} l_k &= \max(h_k^{i-}, h_k^{j-}) - \min(h_k^{i+}, h_k^{j+}) \\ &\leq \max(h_k^i, h_k^j) - \min(h_k^i, h_k^j) = |h_k^i - h_k^j| \\ &\leq \max(h_k^{i+}, h_k^{j+}) - \min(h_k^{i-}, h_k^{j-}) = u_k. \end{aligned} \quad (7)$$

Summing over all histogram bins, we obtain the bounds of the ℓ_1 distance for two unnormalized histograms $\mathbf{h}^i, \mathbf{h}^j$:

$$\sum_k l_k \leq \sum_k |h_k^i - h_k^j| = f_{\ell_1}(\mathbf{h}^i, \mathbf{h}^j) \leq \sum_k u_k. \quad (8)$$

For normalized histograms $\hat{\mathbf{h}}^i, \hat{\mathbf{h}}^j$, we obtain their bounds following same operations of (6) and (7):

$$l_{\ell_1}(\mathbf{R}) = \sum_k \hat{l}_k \leq f_{\ell_1}(\hat{\mathbf{h}}^i, \hat{\mathbf{h}}^j) \leq \sum_k \hat{u}_k = u_{\ell_1}(\mathbf{R}), \quad (9)$$

where

$$\begin{aligned} \hat{l}_k &= \max(\underline{h}_k^i, \underline{h}_k^j) - \min(\overline{h}_k^i, \overline{h}_k^j), \\ \text{and } \hat{u}_k &= \max(\overline{h}_k^i, \overline{h}_k^j) - \min(\underline{h}_k^i, \underline{h}_k^j). \end{aligned} \quad (10)$$

Deriving bounds for ℓ_2 -distance can be written as:

$$\begin{aligned} l_{\ell_2}(\mathbf{R}) &= \sum_k (\hat{l}_k)_+^2 \leq \sum_k (\hat{h}_k^i - \hat{h}_k^j)^2 \\ &= f_{\ell_2}(\hat{\mathbf{h}}^i, \hat{\mathbf{h}}^j) \leq \sum_k \hat{u}_k^2 = u_{\ell_2}(\mathbf{R}), \end{aligned} \quad (11)$$

where $(\cdot)_+ = \max(0, \cdot)$ is a non-negative operator.

2. Histogram intersection Given two normalized histograms, we define their intersection distance by the Hilbert space representation (Scholkopf 2001):

$$f_{\cap}(\hat{\mathbf{h}}^i, \hat{\mathbf{h}}^j) = - \sum_k \min(\hat{h}_k^i, \hat{h}_k^j). \quad (12)$$

Following (5) and (6), we obtain its lower bound and upper bound:

$$\begin{aligned} l_{\cap}(\mathbf{R}) &= - \sum_k \min(\overline{h}_k^i, \overline{h}_k^j) \\ &\leq f_{\cap}(\hat{\mathbf{h}}^i, \hat{\mathbf{h}}^j) \leq - \sum_k \min(\underline{h}_k^i, \underline{h}_k^j) = u_{\cap}(\mathbf{R}). \end{aligned} \quad (13)$$

3. χ^2 distance The χ^2 distance has been proven to be effective to measure distance between histograms. The χ^2 distance is defined as:

$$f_{\chi^2}(\hat{\mathbf{h}}^i, \hat{\mathbf{h}}^j) = \sum_k \frac{(\hat{h}_k^i - \hat{h}_k^j)^2}{\hat{h}_k^i + \hat{h}_k^j}. \quad (14)$$

Incorporating the ℓ_1 -bounds \hat{l}_k and \hat{u}_k in (10) and the inequalities in (5), we obtain the lower bound and upper bound for f_{χ^2} as:

$$l_{\chi^2}(\mathbf{R}) = \sum_k \frac{(\hat{l}_k)_+^2}{\hat{h}_k^i + \hat{h}_k^j}, \quad (15)$$

$$\text{and } u_{\chi^2}(\mathbf{R}) = \sum_k \frac{\hat{u}_k^2}{\hat{h}_k^i + \hat{h}_k^j}. \quad (16)$$

4. Cosine similarity Treating two normalized histograms $\hat{\mathbf{h}}^i$ and $\hat{\mathbf{h}}^j$ as two vectors in the inner product space, we can measure the similarity as their included cosine angle:

$$\begin{aligned} f_C(\hat{\mathbf{h}}^i, \hat{\mathbf{h}}^j) &= \frac{\hat{\mathbf{h}}^i \cdot \hat{\mathbf{h}}^j}{\|\hat{\mathbf{h}}^i\| \|\hat{\mathbf{h}}^j\|} = \frac{\sum_k \frac{h_k^i h_k^j}{|\mathbf{S}^i| |\mathbf{S}^j|}}{\sqrt{\sum_k \left(\frac{h_k^i}{|\mathbf{S}^i|}\right)^2} \sqrt{\sum_k \left(\frac{h_k^j}{|\mathbf{S}^j|}\right)^2}} \\ &= \frac{\sum_k h_k^i h_k^j}{\sqrt{\sum_k (h_k^i)^2} \sqrt{\sum_k (h_k^j)^2}} = \frac{\mathbf{h}^i \cdot \mathbf{h}^j}{\|\mathbf{h}^i\| \|\mathbf{h}^j\|}. \end{aligned} \quad (17)$$

Using (4) and the fact that $\|\mathbf{S}^{i-}\| \leq \|\mathbf{S}^i[b_i, e_i]\| \leq \|\mathbf{S}^{i+}\|$, we obtain the bounds:

$$\begin{aligned} l_C(\mathbf{R}) &= \frac{\sum_k h_k^{i-} h_k^{j-}}{\|\mathbf{S}^{i+}\| \|\mathbf{S}^{j+}\|} \\ &\leq f_C(\mathbf{h}^i, \mathbf{h}^j) \leq \frac{\sum_k h_k^{i+} h_k^{j+}}{\|\mathbf{S}^{i-}\| \|\mathbf{S}^{j-}\|} = u_C(\mathbf{R}). \end{aligned} \quad (18)$$

5. Symmetrized KL divergence By definition, the normalized histograms $\hat{\mathbf{h}}^i$ and $\hat{\mathbf{h}}^j$ are non-negative and sum to one, and thus can be interpreted as two discrete probability distributions. Their similarity can be measured using the symmetrized KL divergence:

$$\begin{aligned} f_D(\hat{\mathbf{h}}^i, \hat{\mathbf{h}}^j) &= D_{KL}(\hat{\mathbf{h}}^i \|\hat{\mathbf{h}}^j) + D_{KL}(\hat{\mathbf{h}}^j \|\hat{\mathbf{h}}^i) \\ &= \sum_k (\hat{h}_k^i - \hat{h}_k^j) (\ln \hat{h}_k^i - \ln \hat{h}_k^j), \end{aligned} \quad (19)$$

where $D_{KL}(\hat{\mathbf{h}}^i \|\hat{\mathbf{h}}^j)$ is the KL divergence of $\hat{\mathbf{h}}^j$ from $\hat{\mathbf{h}}^i$. From (5) and that $\underline{h}_k^i - \overline{h}_k^j \leq \hat{h}_k^i - \hat{h}_k^j \leq \overline{h}_k^i - \underline{h}_k^j$, we have $\ln \underline{h}_k^i - \ln \overline{h}_k^j \leq \ln \hat{h}_k^i - \ln \hat{h}_k^j \leq \ln \overline{h}_k^i - \ln \underline{h}_k^j$. Then, we obtain the bounds for (19):

$$\begin{aligned} l_D(\mathbf{R}) &= \sum_k \left(\underline{h}_k^i - \overline{h}_k^j \right)_+ \left(\ln \underline{h}_k^i - \ln \overline{h}_k^j \right)_+ \\ &\leq f_D(\hat{\mathbf{h}}^i, \hat{\mathbf{h}}^j) \leq \sum_k \left(\overline{h}_k^i - \underline{h}_k^j \right) \left(\ln \overline{h}_k^i - \ln \underline{h}_k^j \right) \\ &= u_D(\mathbf{R}). \end{aligned} \quad (20)$$

6. Symmetrized cross entropy The symmetrized cross entropy (Murphy 2012) measures the average number of bins needed to identify an event by treating each other as the true

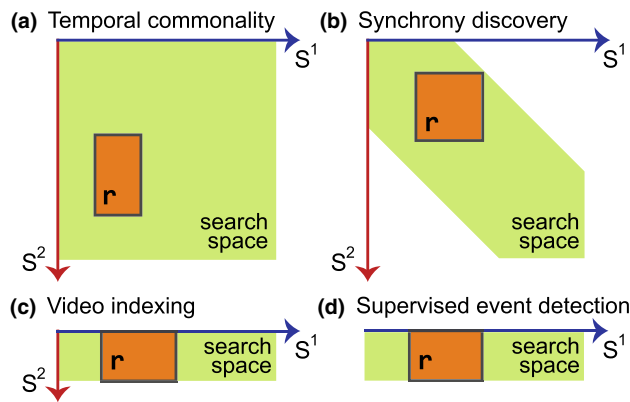


Fig. 3 Searching scenarios readily applicable to the proposed B&B framework: **a** common event discovery (CED), **b** synchrony discovery (SD), **c** video green (VS) and **d** supervised segment-based event detection (ED). Green (lighter) area indicates the search space; an orange (darker) box indicates a candidate solution r . (see Sect. 4 for details) (Color figure online)

distribution. Similar to KL divergence that treats $\hat{\mathbf{h}}^i$ and $\hat{\mathbf{h}}^j$ as two discrete probability distributions, the entropy function is written as:

$$f_E(\hat{\mathbf{h}}^i, \hat{\mathbf{h}}^j) = \sum_k \hat{h}_k^i \log \frac{1}{\hat{h}_k^j} + \sum_k \hat{h}_k^j \log \frac{1}{\hat{h}_k^i}. \quad (21)$$

Recall (5) and that $0 \leq \hat{h}_b^i \leq 1, 0 \leq \hat{h}_b^j \leq 1$, we obtain the bounds:

$$\begin{aligned} l_E(R) &= \sum_b \left(-\hat{h}_b^i \log \hat{h}_b^j - \hat{h}_b^j \log \hat{h}_b^i \right) \\ &\leq f_E(\hat{\mathbf{h}}^i, \hat{\mathbf{h}}^j) \leq \sum_k \left(-\hat{h}_k^i \log \hat{h}_k^j - \hat{h}_k^j \log \hat{h}_k^i \right) \\ &= u_E(R). \end{aligned} \quad (22)$$

Above we have reported derivations for six commonly used measures. However, choice of one or another is influenced by a variety of factors, such as the nature of the data, problem, preferences of individual investigators, etc. In experiments, we picked ℓ_1 , χ^2 , and KL-divergence because due to their popularity in computer vision applications. For instance, ℓ_1 -distance is popular in retrieval problems (e.g., Gusfield 1997; Rubner et al. 2000), χ^2 -distance in object recognition (e.g., Everingham et al. 2006; Lampert et al. 2009), and KL-divergence in measuring similarity between distributions (e.g., Gaussian mixtures for image segmentation Goldberger et al. 2003).

4 Searching Scenarios

With the B&B framework and various bounds derived in the previous section, this section discusses unsupervised and

Algorithm 2: Synchrony Discovery (SD)

Input : A synchronized video pair S^1, S^2 ; minimal discovery length ℓ ; commonality period T
Output: Optimal intervals $r^* = [b_1, e_1, b_2, e_2]$

```

1  $L \leftarrow T + \ell$ ; // The largest possible searching period
2  $Q \leftarrow$  empty priority queue; // Initialize  $Q$ 
3 for  $t \leftarrow 1$  to  $(n - T - L + 1)$  do
4    $R \leftarrow$ 
4      $[t, t+T] \times [t+\ell-1, t+T+L-1] \times [t-T, t+T] \times [t-T+\ell-1, t+T+L-1]$ ;
5    $Q.push(\text{bound}(R), R)$ ;
6 end
7  $r^* \leftarrow \text{BnB}(Q, R)$ ; // BnB procedure in Algo. 1
8 return  $r^*$ ;
```

supervised searching scenarios that can be readily applied. Figure 3 illustrates the searching scenarios in terms of different applications. The first application, common event discovery (CED), as has been discussed in Sect. 3, has the most general form and the broadest search space. Below we discuss others in turn.

4.1 Synchrony Discovery (SD)

Social interaction plays an important and natural role in human behavior. This section presents that a slight modification of CED can result in a solution to discover *interpersonal synchrony*, which is referred as to two or more persons performing common actions in overlapping video frames or segments. Figure 3b illustrates the idea. Specifically, synchrony discovery searches for commonalities (or matched states) among two synchronized videos S^1 and S^2 with n frames each. Rewriting (3), we formulate SD as:

$$\begin{aligned} \text{SD} \quad & \max_{\{b_1, e_1, b_2, e_2\}} f(\varphi_{S^1[b_1, e_1]}, \varphi_{S^2[b_2, e_2]}), \\ & \text{subject to } \ell \leq e_i - b_i, \quad \forall i \in \{1, 2\}, |b_1 - b_2| \leq T, \end{aligned} \quad (23)$$

where $f(\cdot, \cdot)$ is the commonality measure, and T is a *temporal offset* that allows SD to discover commonalities within a T -frame temporal window, e.g., in mother-infant interaction, the infant could start smiling after the mother leads the smile for a few seconds. A naive solution has complexity $\mathcal{O}(n^4)$.

Algorithm: For an event to be considered as a synchrony, they have to occur within a temporal neighborhood between two videos. For this reason, we only need to search within neighboring regions in the temporal search space. Unlike CED or ESS (Lampert et al. 2009) that exhaustively prunes the search space to a unique solution, we constrain the space before the search begins. In specific, we slightly modify Algorithm 1 to solve SD. Let $L = T + \ell$ be the largest possible period to search, we initialize a priority queue Q with rectangle sets $\{[t, t+T] \times [t+\ell-1, t+T+L-1] \times [t-T, t+T] \times [t-T+\ell-1, t+T+L-1]\}_{t=1}^{n-T-L+1}$ and their

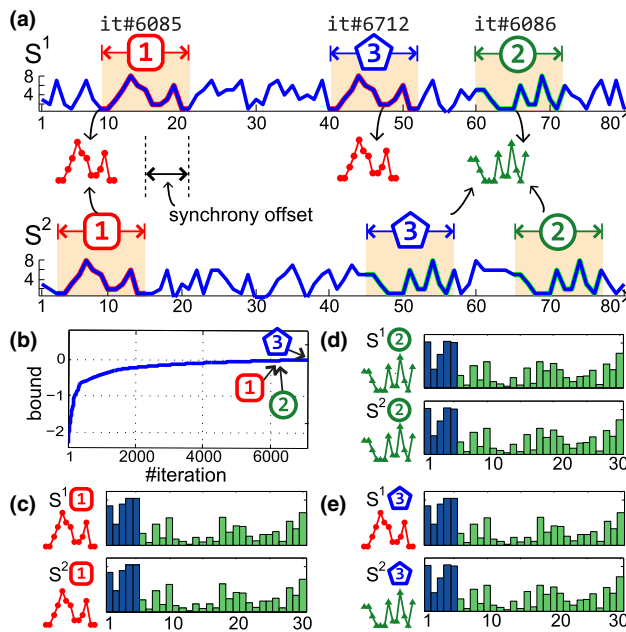


Fig. 4 An example of SD on two 1-D time series using $\ell = 13$ and $T = 5$: **a** top 3 discovered synchronies at different iterations; exhaustive search takes 39,151 iterations. **b** The convergence curve w.r.t. bounding value and #iter. **c–e** Discovered synchronies and their histograms, where blue (darker) and green (lighter) bars indicate the segment features ϕ^{obs} and ϕ^{int} , respectively. ϕ^{int} is $10\times$ magnified for display purpose. The ℓ_1 distances between the three histogram pairs are $6.3\text{e-}8$, $1.5\text{e-}7$, and $5.8\text{e-}2$, respectively (Color figure online)

associated bounds (see details in Sect. 3.4). These rectangle sets lie sparsely along the diagonal in the 2-D search space, and thus prune a large portion before the search. Once all rectangle sets are settled, the CED algorithm can be employed to find the exact optimum. Algorithm 2 summarizes the SD procedure.

Figure 4 shows a synthetic example of 1-D time series with two synchronies, denoted as red dots and green triangle, where one is a random permutation of another. SD discovered 3 dyads with the convergence curve in (b), and histograms of each dyad in (c)–(e). Note that the interaction feature distinguishes the temporal consistency for the first and second discovery, maintaining a much smaller distance than the third discovery.

4.2 Video Search (VS)

The CED algorithm can be also useful for efficient searching for a time series with similar content. That is, given a query time series, search for common temporal segments in a longer video in an efficient manner. Figure 3c illustrates the idea. More formally, let Q be the query time series with length ℓ , we find in the target time series S by modifying (3) as:

Algorithm 3: Video Search (VS)

Input : A query Q with length ℓ ; a target time series S with length n ; a similarity threshold ϵ
Output: Detected events $\{S[b_i, e_i]\}_i$

```

1  $Q \leftarrow$  empty priority queue; // Initialize  $Q$ 
2  $R \leftarrow [1, n] \times [1, n] \times [1, 1] \times [\ell, \ell]$ ; // Initialize  $R$ 
3 while true do
4    $r \leftarrow \text{BnB}(Q, R)$ ; // Obtain  $r$  using BnB (Algo. 1)
5    $b \leftarrow r[0]$ ,  $e \leftarrow r[1]$ ;
6   if  $f(\varphi_{S[b, e]}, \varphi_Q) \leq \epsilon$  then
7     break;
8   end
9   Insert  $S[b, e]$  into  $\{S[b_i, e_i]\}_i$ ;
10   $Q \leftarrow \text{prune}(Q, r)$ ; // Prune space (Sect. 5)
11   $R \leftarrow Q.\text{pop}()$ ;
12 end
13 return  $\{S[b_i, e_i]\}_i$ ;
```

$$\text{VS} \quad \max_{b, e} f(\varphi_{S[b, e]}, \varphi_Q),$$

$$\text{subject to } \ell \leq e - b. \quad (24)$$

The problem now becomes searching along one axis of the search space, but it is still non-convex and non-differentiable. Nevertheless, Algorithm 1 can be directly applied to find the optimal solution by fixing the beginning and ending frame of the query time series, as summarized in Algorithm 3. Note that we do not claim that VS is state-of-the-art method for video search, but just illustrate the versatility of the B&B framework. We refer interested readers to Hu et al. (2011) for a more comprehensive survey.

4.3 Segment-Based Event Detection (ED)

Efficiently detecting variable-length events in time series arises in a wide spectrum of applications, ranging from diseases, financial decline, speech recognition to video security. While event detection has been studied extensively in the literature, little attention has been paid to efficient inference from a pre-trained classifier. Figure 3d illustrates the idea. Here we demonstrate event detection using an SVM decision function, which has been shown effective in many event detection tasks (Schuller and Rigoll 2006; Hoai et al. 2011; Laptev et al. 2008; Sadanand and Corso 2012).

Given the BoTW representation discussed in Sect. 3.1, we represent time series by their histograms. These histograms are used to train an SVM classifier to tell whether a new time series contains an event of interest. To perform inference, temporal segmentation (Schuller and Rigoll 2006; Laptev et al. 2008; Sadanand and Corso 2012) or dynamic programming (DP) (Hoai et al. 2011) is required. However, temporal segmentation for many real-world videos may not be trivial, and DP is computationally expensive to run it in large scale, especially when a time series is too long and relatively small

Algorithm 4: Seg.-based Event Detection (ED)

Input : A video \mathbf{S} of length n ; a pre-trained linear classifier \mathbf{w}
Output: Detected events $\{\mathbf{S}[b_i, e_i]\}_i$

```

1  $Q \leftarrow$  empty priority queue; // Initialize  $Q$ 
2  $R \leftarrow [1, n] \times [1, n]$ ; // Initialize  $R$ 
3 while  $true$  do
4    $r \leftarrow \text{BnB}(Q, R)$ ; // Obtain  $r$  using BnB (Algo. 1)
5    $b \leftarrow r[0], e \leftarrow r[1]$ ;
6   if  $f_{\mathbf{w}}(\mathbf{S}[b, e]) \leq 0$  then
7     break;
8   end
9   Insert  $\mathbf{S}[b, e]$  into  $\{\mathbf{S}[b_i, e_i]\}_i$ ;
10   $Q \leftarrow \text{prune}(Q, r)$ ; // Prune space (Sect. 5)
11   $R \leftarrow Q.\text{pop}()$ ;
12 end
13 return  $\{\mathbf{S}[b_i, e_i]\}_i$ ;

```

portion of frames contain an interested event. Instead, we modify (3) for efficient inference of event detection:

$$\begin{aligned} \text{ED} \quad & \max_{b, e} f_{\mathbf{w}}(\varphi_{\mathbf{S}[b, e]}), \\ & \text{subject to } \ell \leq e - b, \end{aligned} \quad (25)$$

where \mathbf{w} is a pre-trained linear classifier with each element $w_j = \sum_i \alpha_i h_j^i$, and $f_{\mathbf{w}}(\cdot) = \sum_i \alpha_i \langle \cdot, \mathbf{h}^i \rangle$ is the commonality measure based on the classifier. α_i is the weight vector learned during SVM training.

Algorithm The ED problem in (25) becomes supervised detection rather than unsupervised as mentioned in previous sections. The proposed bounds in Sect. 3.4 are thus inapplicable. Due to the summation property of BoTW in (1), we decompose the commonality measure into per-frame positive and negative contributions: $f_{\mathbf{w}}(\mathbf{S}[b, e]) = \sum_{i=b}^e (f_{\mathbf{w}}^+(\mathbf{S}[i, i]) + f_{\mathbf{w}}^-(\mathbf{S}[i, i]))$. Denote the longest and the shortest possible searching segments as \mathbf{S}^+ and \mathbf{S}^- respectively, with slight abuse of notation, we reach the bounds:

$$\begin{aligned} l_{\mathbf{w}}(R) &= f_{\mathbf{w}}^+(\mathbf{S}^-) + f_{\mathbf{w}}^-(\mathbf{S}^+) \\ &\leq f_{\mathbf{w}}^+(\mathbf{S}) + f_{\mathbf{w}}^-(\mathbf{S}) = f_{\mathbf{w}}(\mathbf{S}) \\ &\leq f_{\mathbf{w}}^+(\mathbf{S}^+) + f_{\mathbf{w}}^-(\mathbf{S}^-) = u_{\mathbf{w}}(R), \end{aligned} \quad (26)$$

where $R = [b, e]$ corresponds to time series \mathbf{S} , instead of previous definition over two time series. With the derived bounds, the CED algorithm can be directly applied for efficient inference of an event of interest, as summarized in Algorithm 4.

4.4 Comparisons with Related Work

The proposed CED bear similarities and differences with several related work. Below we discuss in terms of problem definition and technical details.

Problem definition Although CED achieves discovery via “matching” between subsequences, it has fundamental differences from standard matching problems. For instance, CED allows *many-to-many* mapping (e.g., Sect. 6.1.2), while standard matching algorithms assume *one-to-one* or *one-to-many* mapping. Moreover, a matching problem (e.g., graph matching or linear assignment) typically measures sample-wise similarity or distance to determine correspondence between one another, e.g., a feature vector on a node in a graph. CED uses bag-of-words representation that aggregates multiple samples (i.e., frames) into one vector, making the application of standard matching methods non-trivial.

CED is also different from time warping (e.g., dynamic time warping Keogh and Ratanamahatana 2005) and temporal clustering (e.g., aligned cluster analysis Zhou et al. 2013). Time warping aims to find the optimal match between two given sequences that allow for stretched and compressed sections of the sequences. Given this goal, time warping assumes the beginning and the ending frames of the sequences to be fixed, and performs matching on entire sequence. Similarly, temporal clustering considers entire sequence in its objective, and hence is likely to include irrelevant temporal segments in one cluster. On the contrary, CED does not assume fixed beginning and ending frames, instead directly targeting at subsequence-subsequence matching, and thus enables a large portion of irrelevant information to be ignored.

Technical details Technically, the proposed B&B framework is closely related to Efficient Subwindow Search (ESS) (Lampert et al. 2009) and Spatio-Temporal B&B (STBB) (Yuan et al. 2011). However, they have at least three differences. (1) *Learning framework*: ESS and STBB are supervised techniques that seek for a confident region in image or a volume in video according to a pre-trained classifier. CED is unsupervised, and thus requires no prior knowledge. (2) *Bounding functions*: We design new bounding functions for the unsupervised CED problem. Moreover, ESS and STBB consider only upper bounds, while CED can incorporate both upper and lower bounds. (3) *Search space*: ESS and STBB search over spatial coordinates of an image or a spatio-temporal volume in a video, while CED focuses on temporal positions over time series.

For segment-based event detection (ED), we acknowledge its similarity with the version of STBB that omits spatial volume. Both address efficient search in a one-dimension time series, and differ in the following ways. (1) *Objective*: ED searches for segments with maximal, positive *segment-based* decision values. STBB uses a Kadane’s algorithm for *frame-based* max subvector search, which potentially lead to inferior detection performance because the max sum is usually found in an overly-large segment (as can be seen in Sect. 6.3). (2) *Searching strategy*: ED prunes the search space to avoid evaluating segments where an AU is unlikely to occur; STBB evaluates every frame. (3) *Inputs*: ED can

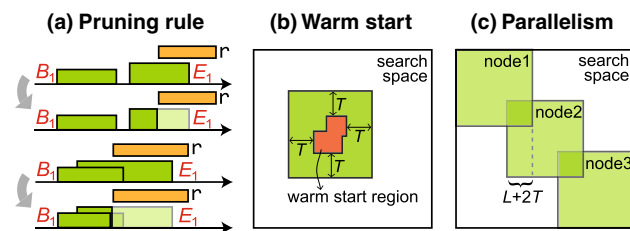


Fig. 5 Illustration of extensions: **a** pruning rules applied to multiple-commonality discovery, **b** SD with warm start, and **c** SD with parallelism

take the minimal length and normalized histograms as input, yet it is unclear for STBB to accommodate such input because of the linear nature of the Kadane's algorithm.

5 Extensions to the B&B Framework

Given the aforementioned CED algorithm and variants, this section describes extensions to discovery among multiple time series and discover multiple commonalities. Due to the special diagonal nature of SD, we also introduce its acceleration using warm start and parallelism. Figure 5 illustrates these extensions.

Discovery among multiple time series We have described above how the B&B framework can discover temporal commonalities within a pair of time series. Here we show that the framework can be directly extended to capture commonality among multiple time series. Specifically, we formulate the discovery among N sequences $\{S^i\}_{i=1}^N$ by rewriting (3) as:

$$\begin{aligned} \max_{\{b_i, e_i\}_{i=1}^N} & F\left(\{\phi_{S^i[b_i, e_i]}\}_{i=1}^N\right) \\ \text{subject to } & \ell \leq e_i - b_i, \quad \forall i \in \{1, \dots, N\}, \end{aligned} \quad (27)$$

where $F(\cdot)$ is a similarity measure for a set of sequences and defined as the sum of pairwise similarities:

$$F\left(\{\phi_{S^i[b_i, e_i]}\}_{i=1}^N\right) = \sum_{i \neq j} f\left(\phi_{S^i[b_i, e_i]}, \phi_{S^j[b_j, e_j]}\right). \quad (28)$$

Given a rectangle set R and a time series pair (S^i, S^j) , we rewrite their pairwise bounds in Sect. 3.4 as $l_f^{ij}(R)$ and $u_f^{ij}(R)$. The bounds for $F(\cdot, \cdot)$ can be defined as:

$$\begin{aligned} l_F(R) &= \sum_{i \neq j} l_f^{ij}(R) \leq F\left(\{\phi_{S^i[b_i, e_i]}\}_{i=1}^N\right) \\ &\leq \sum_{i \neq j} u_f^{ij}(R) = u_F(R). \end{aligned} \quad (29)$$

Given this bound, Algos. 1 and 2 can be directly applied to discover commonalities among multiple time series.

Discover multiple commonalities Multiple commonalities occur frequently in real videos, while the B&B framework only outputs one commonality at a time. Here, we introduce a strategy that prunes the search space to accelerate multiple commonality discovery. Specifically, we repeat the searching algorithm by passing the priority queue Q from the previous search to the next, and continue the process until a desired number of solutions is reached, or the returned commonality measure $f(\cdot, \cdot)$ is less than some threshold. The threshold can be also used for excluding undesired discoveries for the scenario where two sequences have no events in common. That is, if the first discovery does not pass a pre-defined threshold, the algorithm returns empty because the subsequent discoveries perform no better than the first one. Figure 5a illustrates an example of the pruning rule when E_1 overlaps with a previously discovered solution r . Because we want to exclude the same solution for the next discovery, the search region is updated by avoiding overlapping with previous solution. For axes of both S^1 and S^2 , all R overlapped with r is updated using the same rule, or discarded if the updated R is empty, i.e., $|R| = 0$. The updated rectangle sets, along with their bounds, are then pushed back to Q before the next search.

This pruning strategy is simple yet very effective. Previously derived bounds remain valid because each updated set is a subset of R . In practice, it dramatically reduces $|Q|$ for searching the next commonality. For example, in synchrony discovery of Fig. 4, $|Q|$ is reduced 19% for the second search, and 25% for the third SD. Note that this pruning strategy differs from conventional detection tasks, e.g., Lampert et al. (2009), Yuan et al. (2011), which remove the whole spatial or temporal region for the next search. In CED, temporal segments can be many-to-many matching, i.e., $S^1[b_1, e_1]$ can match multiple segments in S^2 and vice versa. Thus, removing any segments from either time series would cause missing matches. This strategy allows us to maintain *many-to-many* matching.

SD with Warm start Due to the B&B nature, SD exhibits poor worst-case behavior, leading to a complexity as high as an exhaustive search (Narendra and Fukunaga 1977). On the other hand, B&B can quickly identify the exact solution when a local neighborhood contains a clear optimum (Lampert et al. 2009). Given this motivation, we explore a “warm start” strategy that estimates an initial solution with high quality, and then initializes SD around the solution. Estimating an initial solution costs only few percentage of total iterations, and thus can effectively prune branches in the main SD algorithm. Figure 5b illustrates the idea. Specifically, we run sliding window sampled with step-size=10, sort the visited windows according their distances, and then determine a warm start region around the windows within the top one percentile. Then SD is performed only within an expanded neighborhood around the warm start region.

Table 1 Distribution of event lengths in different datasets: min and max show the shortest and longest length of a common event

Dataset	Min	25-th	50-th	75-th	Max	Std
RU-FACS (Bartlett et al. 2006)	13	42	79	159	754	125.6
Mocap (http://mocap.cs.cmu.edu/)	41	142	175	218	483	67.5

25-, 50-, and 75-th indicate degrees of percentiles

SD with Parallelism The use of parallelism to speed up B&B algorithms has emerged as a way for large problems (Gendron and Crainic 1994). Based on the block-diagonal structure in the SD search space, this section describes an parallelized approach to scale up SD for longer time series. In specific, we divide SD into subproblems, and perform the SD algorithm solve each in parallel. Because each subproblem is smaller than the original one, the number of required iterations can be potentially reduced. As illustrated in Fig. 5c, the original search space is divided into overlapping regions, where each can be solved using independent jobs on a cluster. The results are obtained as the top k rectangles collected from each subproblem. Due to the diagonal nature of SD in the search space, the final result is guaranteed to be a global solution. The proposed structure enables static overload distribution, leading to an easily programmable and efficient algorithm.

6 Experiments

In this section, we evaluated the effectiveness and efficiency of the proposed B&B framework under three applications: Common event discovery (Sect. 6.1), synchrony discovery (Sect. 6.2), and variable-length segment-based event detection (Sect. 6.3). As mentioned in Sect. 4, each application relates to a particular searching scenario of the B&B framework.

6.1 Common Event Discovery (CED)

In the first experiment, we evaluated CED on discovering common facial events, and discovering multiple common human actions.

Table 1 shows the distribution of event lengths in respective experiments. The mixture of long and short events indicates a more realistic scenario of handling events with slow and fast motions. Specifically, for RU-FACS, we computed the distribution of AU12 events among the 4950 sequence pairs. For mocap, the distribution was computed on a total of 25 actions from 45 sequence pairs (details below).

6.1.1 Discovering Common Facial Events

This experiment evaluates the CED algorithm to find similar facial events in the RU-FACS dataset (Bartlett et al. 2006).

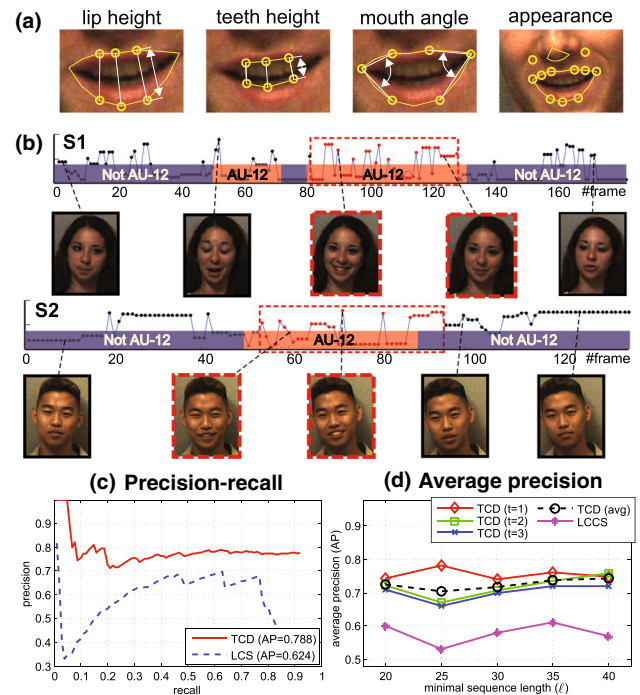


Fig. 6 Results on discovering common facial actions: **a** Facial features extracted from the tracked points. **b** An example of common discovered facial events (indicated by dashed-line rectangles). **c**, **d** Accuracy evaluation on precision-recall and average precision (AP)

The RU-FACS dataset consists of digitized video of 34 young adults. They were recorded during an interview of approximately 2 min duration in which they lied or told the truth in response to interviewer's questions. Pose orientation was mostly frontal with moderate out-of-plane head motions. We selected the annotation of Action Unit (AU) 12 (*i.e.*, mouth corner puller) from 15 subjects that had the most AU occurrence. We collected 100 video segments containing one AU 12 and other AUs, resulting in 4950 pairs of video clips from different subjects. For each video, we represented features as the distances between the height of lips and teeth, angles for the mouth corners and SIFT descriptors in the points tracked with Active Appearance Models (AAM) (Matthews and Baker 2004) (see Fig. 6a for an illustration).

Accuracy evaluation Because the CED problem is relatively new in computer vision, to our knowledge there is no baseline we could directly compare to. Instead, we compared against the state-of-the-art sequence matching approach: Longest common consecutive subsequence matching (LCCS) (Wang and Velipasalar 2009). Observe that when

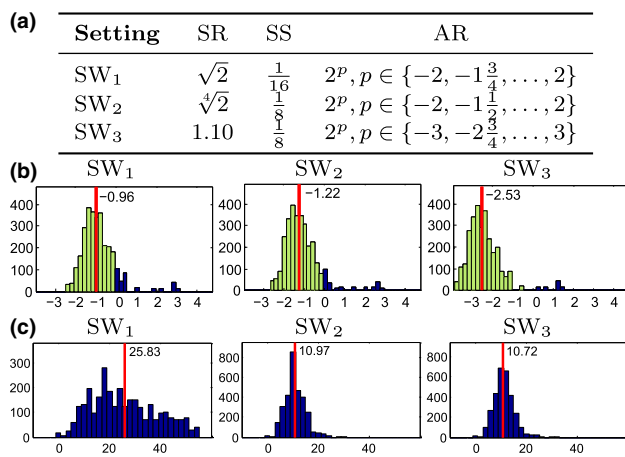


Fig. 7 Efficiency evaluation between CED and alternative sliding window (SW) approach. **a** Parameter settings (Viola and Jones 2004; Lampert et al. 2009): size-ratio (SR), stepsize (SS), and aspect ratios (AR). **b** Histogram of ratio of #evaluation: $\log \frac{n_{\text{CED}}}{n_{\text{SW}_i}}$. Red vertical lines indicate the average. Light green bars show CED performs less evaluations than SW; dark blue bars represent the opposite. **c** Histogram of differences between resulting commonality measure: $f_{\ell_1}(\mathbf{x}^{\text{SW}_i}) - f_{\ell_1}(\mathbf{x}^{\text{CED}})$ (Color figure online)

the per-frame feature was quantized into a temporal word, the unsupervised CED problem can be naturally interpreted as an LCCS. Following LCCS that uses a 0–1 distance, we chose ℓ_1 -distance for CED. Note that the segment-based BoTW representation is not helpful for LCCS (Wang and Velipasalar 2009), because LCCS computes matches only at frame-level. The minimal length ℓ was fixed as the smaller length of ground truth segments for both LCCS and CED. Given a discovered solution \mathbf{x} and a ground truth \mathbf{g} that indicates a correct matching, we measured their *overlap score* (Everingham et al. 2006) as $\text{overlap}(\mathbf{x}, \mathbf{g}) = \frac{\text{area}(\mathbf{x} \cap \mathbf{g})}{\text{area}(\mathbf{x} \cup \mathbf{g})}$. The higher the overlap score, the better the algorithm discovered the commonality. We considered \mathbf{x} to be a correct discovery if the overlap score is greater than 0.5.

Figure 6b shows an example of a correct discovery of AU12. In this example, CED was able to correctly locate an AU 12 segment with overlap score greater than 0.8. Figure 6c plots the precision-recall curves for the first discovery of CED and LCCS. We reported the average precision (AP) (Everingham et al. 2006) and found CED outperformed LCCS by 0.15 points. Unlike LCCS that sought for identical subsequences, CED considered a distribution of temporal words present in two videos, and thus was able to more reliably capture common events in real-world videos. Figure 6d shows the average precision of our approach under different parameters. We varied the minimal sequence length ℓ in $\{20, 25, \dots, 40\}$, and examined the AP of the t -th result. As can be observed from the averaged AP (black dashed line), our B&B approach performed more stably across different combinations of ℓ and t . As a result, CED performed on average 16% higher AP than LCCS in discovering the common facial actions.

Efficiency evaluation Using the above settings, we evaluated speedup of the CED algorithm against exhaustive sliding window (SW) approach, which was implemented following parameter settings in Lampert et al. (2009), Viola and Jones (2004). Figure 7a shows these settings denoted as SW_{*i*} ($i = 1, 2, 3$). Denote lengths of two time series as m, n and the minimal length for each sequence is ℓ , we set the maximal and minimal rectangle size for SW to be $(m \times n)$ and $(\ell \sqrt{\text{AR}} \times \frac{\ell}{\sqrt{\text{AR}}})$, respectively. To be independent of implementation, we measured the *discovery speed* as the number of evaluation for the bounding functions, referred as n^{CED} and n^{SW_i} for CED and SW_{*i*} respectively. Figure 7b shows the histograms of the log ratio for $n^{\text{CED}}/n^{\text{SW}_i}$. The smaller the value, the less times CED has to evaluate the distance function. As can be seen, although SW was parameterized to search only a subset of the search space, CED searched the entire space yet still performed on average 6.18 times less evaluations than SW. To evaluate the *discovery quality*, we computed the distance difference measured by CED and SW, i.e., $f_{\ell_1}(\mathbf{x}^{\text{SW}_i}) - f_{\ell_1}(\mathbf{x}^{\text{CED}})$. The larger the difference, the lower quality of discovery SW got. Figure 7c shows the histograms of such differences. One can observe that the differences are always greater than or equal to zero. This is because our method provably finds the global optimum. On the other hand, SW only performed a partial search according to its parameters, and thus was likely to reach larger distance than ours.

6.1.2 Discover Multiple Common Human Motions

This experiment attempts to discover *multiple* common actions using the CMU-Mocap dataset (<http://mocap.cs.cmu.edu/>). We used Subject 86 that contains 14 long sequences with 1,200~2,600 frames and human action annotation (Barbič et al. 2004). Each sequence contains up to 10 actions (out of a total of 25) such as *walk*, *jump*, *punch*, etc. See Fig. 8a for an example. Each action ranged from 100 to 300 frames. We randomly selected 45 pairs of sequences and discovered common actions among each pair. Each action was represented by root position, orientation and relative joint angles, resulting in a 30-D feature vector. Note that this experiment is much more challenging than the previous one due to the large number of frames and more complicated actions. In this case, we excluded SW for comparison because it needs 10^{12} evaluations that is impractical.

Figure 8a illustrates the first six common motions discovered by CED. A failure discovery is shown in the shaded number 6, which matches *walk* to *kick*. An explanation is because these actions were visually similar, resulting in similar features of joint angles. Figure 8b shows the precision-recall curve for different values of overlapping threshold ε . Using ℓ_1 distance, the curve decreases about 10% AP when the overlap score ε raises from 0.4 to 0.7, which implies that

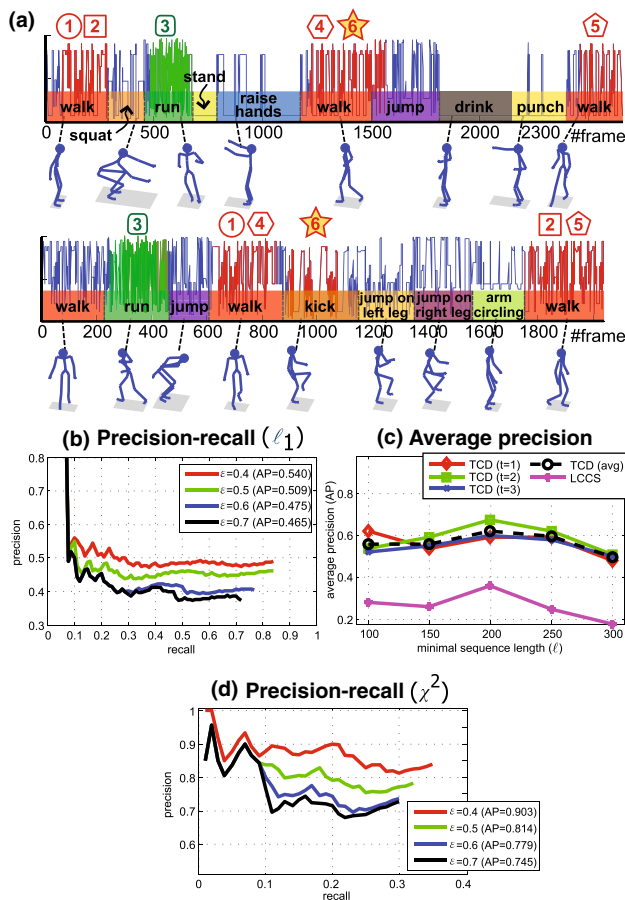


Fig. 8 **a** Top six discovered common motions. The numbers indicate discovered commonalities. Note that the shaded star (number 6) indicates an incorrect discovery that matched walk and kick. **b**, **c** Precision-recall and average precision on ℓ_1 distance. **d** Precision-recall on χ^2 distance

we can retain higher quality results without losing too much precision. Figure 8c shows the average precision over various ℓ on the t -th discovered result. LCCS performed poorly to obtain long common subsequences because human motions have more variability than just one facial event (e.g., AU-12). On the contrary, CED used BoTW representation, and thus allowed more descriptive power for activity recognition. Figure 8d shows the precision-recall curve evaluated with χ^2 distance. Although the Mocap dataset is very challenging in terms of various motions and diverse sequence lengths, the CED algorithm with χ^2 performed 30% better than ℓ_1 and LCCS. It suggests χ^2 is a more powerful commonality measure for histograms than ℓ_1 . Overall, using the χ^2 measurement and $\varepsilon = 0.5$, CED achieved 81% precision.

6.2 Synchrony Discovery (SD)

This section evaluates SD for discovering synchronous behavior using three datasets of increasing diversity: Posed

Table 2 Distance and quality analysis on CMU Mocap dataset: (top) χ^2 distance using $1e-3$ as unit, (bottom) recurrent consistency

Pair	(1,1)	(2,4)	(3,13)	(5,7)	(6,8)	(9,10)	(12,14)	Avg.
χ^2 -distance								
SD	6.3	1.2	4.7	2.6	0.1	0.2	11.9	3.9
SW ₅ [*]	6.5	1.3	6.7	5.4	0.1	0.4	12.0	4.6
SW ₁₀ [*]	6.7	2.7	6.7	10.1	0.2	0.7	14.3	5.9
SW ₅ ^{μ}	97.1	76.9	81.4	64.2	89.3	172.0	334.5	130.8
SW ₅ ^{σ}	33.8	74.4	53.8	28.2	79.2	117.7	345.1	104.6
SW ₁₀ ^{μ}	94.8	77.3	81.8	63.2	87.1	170.2	327.2	128.8
SW ₁₀ ^{σ}	34.3	74.1	54.2	28.3	79.4	117.8	341.5	104.2
Recurrence consistency								
SD	0.89	0.85	0.46	0.90	1.00	0.64	0.76	0.79
SW ₅ [*]	0.95	0.81	0.50	0.84	1.00	0.69	0.73	0.79
SW ₁₀ [*]	0.95	0.75	0.50	0.64	1.00	0.55	0.00	0.63
SW ₅ ^{μ}	0.07	0.32	0.09	0.07	0.08	0.13	0.12	0.12
SW ₅ ^{σ}	0.16	0.33	0.25	0.20	0.21	0.29	0.22	0.24
SW ₁₀ ^{μ}	0.08	0.31	0.09	0.07	0.09	0.13	0.12	0.13
SW ₁₀ ^{σ}	0.19	0.33	0.26	0.21	0.22	0.29	0.23	0.25

SW₅^{*} indicates the optimal window found by SW_s with step size $s = 5, 10$; SW_s ^{μ} and SW_s ^{σ} indicate average and standard deviation among all windows. The best discovery are marked in bold

motion capture (Sec. 6.2.1) and unposed, spontaneous video of mothers and their infants (Sec. 6.2.2) and of young adults in a small social group (Sec. 6.2.3).

6.2.1 Human Actions

We first provide an objective evaluation the SD algorithm (Sec. 6.1.2) on discovering human actions using the CMU Mocap dataset (<http://mocap.cs.cmu.edu/>). Mocap data provides high-degree reliability in measurement and serves as an ideal target for a clean-cut test of our method. To mimic a scenario for SD, we grouped the sequences into 7 pairs as the ones containing similar number of actions, and trimmed each action to up to 200 frames. SD was performed using $\ell = 120$ and $T = 50$. Denote the video index set as \mathcal{A} , we evaluated the discovery performance by the *recurrence consistency* (Delaherche et al. 2012):

$$\mathcal{Q}(\tau) = \frac{1}{C \prod_i n_i} \sum_c \sum_{(i,j) \in \mathcal{A}} \sum_{p,q} I(\mathbf{Y}_i^c[p] = \mathbf{Y}_j^c[q]), \quad (30)$$

where $I(X)$ is an indicator function returning 1 if the statement X is true and 0 otherwise, and $\mathbf{Y}_i^c[p]$ denote the c -th class annotation corresponding to the p -th frame in \mathbf{S}^i .

Table 2 summarizes the SD results compared with the baseline sliding window (SW). Results are reported using χ^2 -distance and the recurrent consistency. A threshold of 0.012 was manually set to discard discovery with large dis-

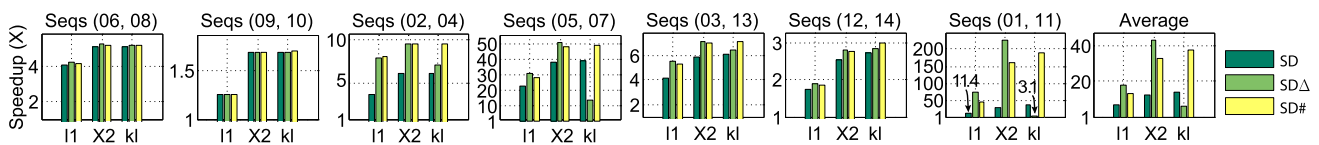


Fig. 9 Speedup of SD against sliding window (SW) on CMU-Mocap. All 7 pairs of sequences from subject 86 were evaluated. The speedup was computed as the relative number of evaluations N^{SW}/N^{SD} using ℓ_1 , χ^2 and symmetrized KL divergence

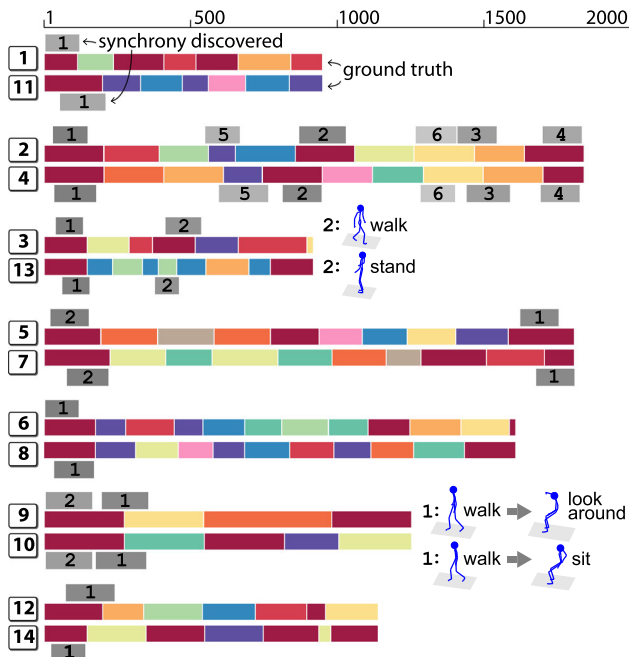


Fig. 10 Discovered synchronies on 7 pairs of Subject 86 in CMU-Mocap dataset. Each pair is annotated with ground truth (colorful bars, each represents an action), and synchronies discovered by our method (shaded numbers). Synchronies with disagreed action labels are visualized (Color figure online)

tance. We ran SW with step sizes 5 and 10, and marked the windows with the minimal distance as SW_5^* and SW_{10}^* , respectively. Among all, SD discovers all results found by SW. To understand how well a prediction by chance can be, all windows were collected to report average μ and standard deviation σ . As can be seen, on average, a randomly selected synchrony can result in large distance over 100 and low quality below 0.3. SD maintained an exact minimal distance with good qualities as the ones found by exhaustive SW. Note that, because SD is totally unsupervised, the synchrony with minimal distance may not necessarily guarantee the highest quality.

Figure 9 shows the speed up of SD against exhaustive SW. SD and its extensions demonstrated an improved efficiency over SW. In some cases, SD^Δ improved search speed by a large margin, *e.g.*, in (01,11) with χ^2 -distance reached a speed boost over 200 times. Across all metrics, the speed up of SD^Δ was less obvious with symmetrized KL divergence. $SD^\#$ was implemented on a 4-core machine; an extension to larger clusters is possible yet beyond the scope of this study.

On average, $SD^\#$ consistently accelerated the original SD due to parallelism.

Figure 10 shows the qualitative results on all 7 pairs, annotated with ground truth and the discovered synchronies. As can be seen, SD allows to discover multiple synchronies with varying lengths. Although some discovered synchronies contain disagreed action labels, one can observe that the discoveries share reasonable visual similarity, *e.g.*, in pair (9,10), the “look around” action in sequence 9 was performed when the subject was seated, sharing the similarity with the “sit” action in sequence 10.

6.2.2 Parent-Infant Interaction

Parent-infant interaction is critical for early social development. This section attempts to characterize their affective engagement by exploring the moments where the behavior of both the parent and the infant are correlated. We performed this experiment on the mother-infant interaction dataset (Messinger et al. 2009). Participants were 6 ethnically diverse 6-month-old infants and their parents (5 mothers, 1 father). Infants were positioned in an infant-seat facing their parent who was seated in front of them. We used 3 min of normal interaction where the parent plays with the infant as they might do at home. Because this dataset was not fully annotated, we only evaluated the results quantitatively. After the faces were tracked, we used only the shape features because the appearance of adults and infants are different. Throughout this experiment, we set $\ell = 80$ and $T = 40$.

Figure 11 illustrates three discovered synchronies among all parent-infant pairs. As can be seen, many synchronies were discovered as the moments when both infants and parents exhibit strong smiles, serving as a building block of early interaction (Messinger et al. 2009). Besides smiles, a few synchronies showed strong engagement in their mutual attention, such as the second synchrony of group ① where the infant cried after the mother showed a sad face, and the second synchrony of the second group where the mother stuck her tongue out after the infant did so. These interactive patterns offered solid evidence of a positive association between infants and their parents.

6.2.3 Social Group Interaction

This experiment investigates discovery of synchronies in social group interaction. We used the GFT dataset (Sayette



Fig. 11 Discovered synchronies from 6 pairs of parents and infants interacting. Each *column* indicates a discovery and its #frame

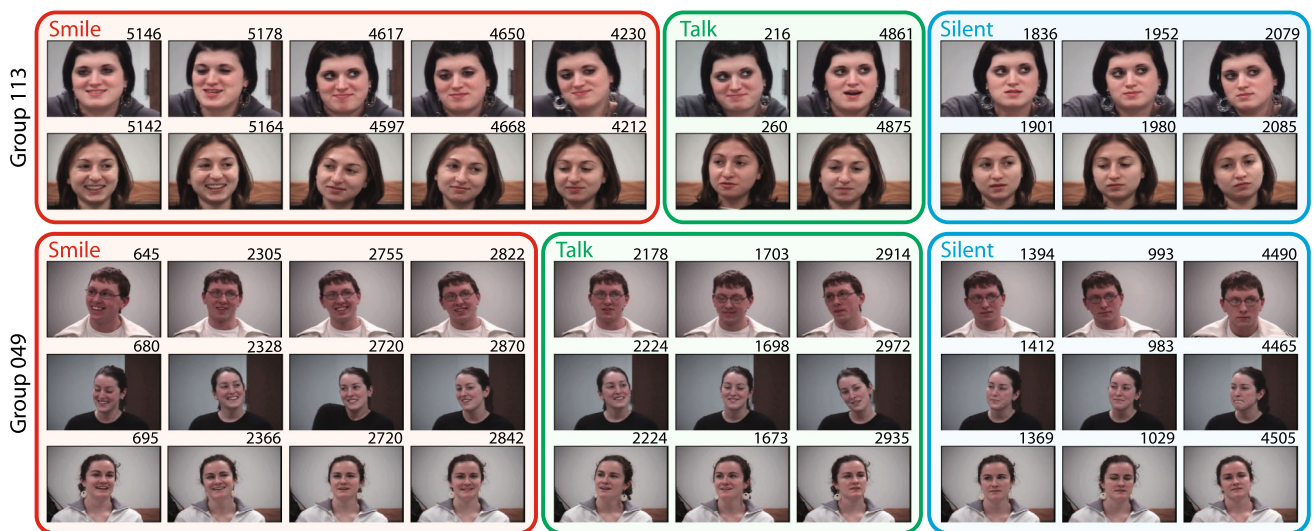


Fig. 12 Top 10 discovered synchronies from groups 113 and 128 in the GFT dataset. Each *column* indicates a discovered synchrony and its frame number. The SD algorithm correctly matched the states of *smiling*, *talking* and *silent*

et al. 2012) that consists of 720 participants recorded during group-formation tasks. Previously unacquainted participants sat together in groups of 3 at a round table for 30 min while getting to know each other. We used 2 min of videos from 48 participants, containing 6 groups of two subjects and 12 groups of three subjects. SD was performed to discover *dyads* among groups of two, and *triads* among groups of three. Each video was tracked with 49 facial landmarks using IntraFace (De la Torre et al. 2015). We represented each face by concatenating appearance features (SIFT) and shape features (49 landmarks). In this dataset, we used annotations of AUs (10,12,14,15,17,23,24) that appear most frequently.

Figure 12 shows qualitative results of the discovered dyadic and triadic synchronies among two social groups. Each column indicates a discovery among each group. As can be observed, most common events are discovered as concurrent smiles, talk, or silent moments where all partici-

pants remained neutral. Because the interaction was recorded during a drinking section, the SD algorithm discovers more frequent concurring smiles than other behavior. This discovery is particular interesting for complying with the findings in Sayette et al. (2012) that alcohol facilitates bonding during group formation. It is noticeable that the SD algorithm requires no human supervision, yet can identify meaningful patterns (e.g., smiles) occult to supervised approaches.

Quantitatively, we examined SD with varying ℓ , i.e., $\ell \in \{30, 60, 120\}$, resulting in synchronies that last at least 1, 2 and 4 s; we set the synchrony offset $T = 30$ (1 s). Baseline SW was performed using step sizes 5 and 10. Symmetrized KL divergence was used as the distance function. We evaluated the distance and quality among the optimal window discovered, as well as the average and standard deviation among all windows to tell a discovery by chance. Figure 13 shows the averaged KL divergence and recurrent consistency

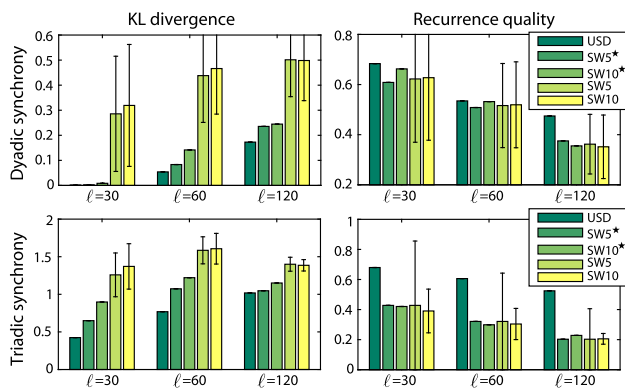


Fig. 13 Analysis on top 10 discovered dyadic and triadic synchronies of the GFT dataset. SW denoted with *asterisk* indicates the optimal windows discovered, and without *asterisk* indicates the average and standard deviation over all visited windows

(Eq. (30)) among top 10 discovered dyadic and triadic synchronies. As can be seen, SD always guarantees the lowest divergence because of its nature to find the exact optimum. The recurrence quality decreases while ℓ grows, showing that finding a synchrony with longer period while maintaining good quality is harder than finding one with shorter period. Note that, although the discover quality is not guaranteed in an unsupervised discovery, SD consistently maintained the best discovery quality across various lengths. This result illustrates the power of our unsupervised method that agrees with that of supervised labels.

6.3 Segment-Based Event Detection (ED)

This experiment evaluates performance and computation time of segment-based event detection on the GFT dataset (Sayette et al. 2012), as used in Sect. 6.2.3. The task is to localize AU events using a pre-trained segment-based linear SVM classifier. The AUs of interest are 1, 2, 6, 7, 10, 11, 12, 14, 15, 17, 23, and 24. Unlike previous studies that require temporal segmentation (Schuller and Rigoll 2006; Laptev et al. 2008; Sadeanand and Corso 2012), we focused on joint detection and segmentation of a temporal event. Specifically, we compared ED with a hybrid SVM-HMM (Krüger et al. 2005) (denoted HMM hereafter for simplicity) and the state-of-the-art event detection algorithms, including a dynamic programming (DP) approach (Hoai et al. 2011) and the Kadane's algorithm used in STBB (Yuan et al. 2011). We trained a frame-based SVM for each AU, and used the same SVM for the detection task on different methods. For SVM-HMM, the HMM has two states, *i.e.*, activation or inactivation of an AU. The state transition probabilities and the a-priori probability were estimated by the frequency of an AU activation in the training data. The emission probabilities of HMM was computed based on normalized SVM output using Platt's scaling (Platt 1999). During test, the most likely AU state path

for each video was determined by a standard Viterbi algorithm, which has a complexity $\mathcal{O}(|s|^2 \times N)$, where $|s| = 2$ is the number of states and N is the number of frames of a test video. For both ED and DP, we set the minimal discovery length $\ell = 30$. For DP, we set the maximal segment lengths in $\{100, 150, 200\}$, denoted as DP_100, DP_150, and DP_200, respectively. For evaluation, we used the standard F1 score and the F1-event metric (Ding et al. 2012) defined as $F1\text{-event} = \frac{2EP \cdot ER}{EP + ER}$, where EP and ER stand for event-based precision and event-based recall. Unlike a standard F1 score, F1-event focuses on capturing the temporal consistency of prediction. An event-level agreement holds if the overlap of two temporal segments is above a certain threshold.

Figure 14a shows the F1-event curve w.r.t. event overlapping thresholds. Overall DP and ED performed better than the baseline HMM. The performance of DP dropped when threshold was greater than 0.6, which implies DP missed highly overlapped events during detection. This is because DP performed exhaustive search, and thus requested a maximal search length for computational feasibility. On the other hand, ED by construction excludes such limitation. Figure 14b shows the running time on a 2.8 GHz dual core CPU machine by comparing ED v.s. DP. Note that we omitted STBB and HMM in Fig. 14b because the time difference between ED and STBB/HMM is insignificant under this scale. Each detected AU event is plotted in terms of the running time and sampled video length (#frame). As can be seen, the computation time for DP increased linearly with video length, while ED maintained invariance of video length. These results suggest that ED was able to perform comparably with significantly improved efficiency for event detection.

Figure 14c, d shows the trend of running time v.s. F1-event and F1 score across ED and all alternative methods. Each marker indicates a detection result for a sequence. For visualization purpose, we randomly picked 120 sequences to include in this figure. The quantitative evaluation on the entire dataset is shown in Table 3. As can be seen in Fig. 14c, d, STBB and HMM performed significantly faster than others due to their linear nature in computation. In general, for F1-event and F1, STBB led to suboptimal performance because events with activation are usually found in over-length segments. Figure 14e illustrates detection results of three subjects. In all cases, it reveals the over-length detection of STBB due to its consideration of max subvectors. As can be seen, STBB tends to include a large temporal window so that the sum of decision values is maximized. HMM took SVM outputs as emission probability, and thus performs similarly as a frame-based SVM. HMM tends to generate lower F1-event, as also suggested in Fig. 14a. This is because of the memoryless property considered in the Markov chain, *i.e.*, the future state only depends upon the present state. On the

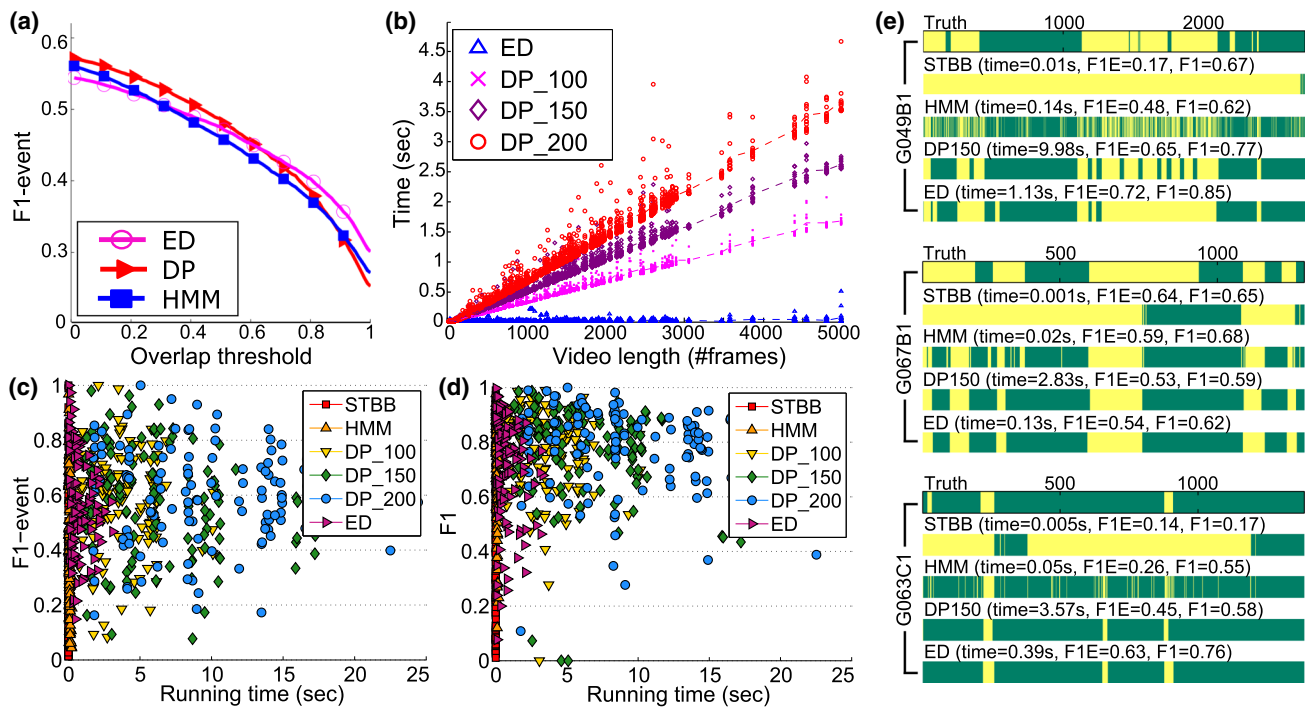


Fig. 14 Comparison between ED and alternative approaches in terms of: **a** F1-event over 12 AUs, **b** running time v.s. video length, **c** F1-event v.s. time, **d** F1 v.s. time and **e** comparison between ground truth

and detection results on 3 subjects. *Light yellow* and *dark green* indicate activation and deactivation of AU12, respectively

Table 3 Comparison between ED and alternative methods in terms of running time, F1-event (F1E), and F1 on supervised AU detection

Method	Time (s)	F1E	F1
STBB	0.003 ± 0.002	0.297 ± 0.256	0.420 ± 0.270
HMM	0.090 ± 0.049	0.405 ± 0.209	0.698 ± 0.182
DP100	3.987 ± 2.184	0.586 ± 0.188	0.756 ± 0.179
DP150	6.907 ± 3.720	0.586 ± 0.188	0.756 ± 0.179
DP200	9.332 ± 5.268	0.586 ± 0.188	0.756 ± 0.179
ED (ours)	0.668 ± 0.873	0.572 ± 0.197	0.753 ± 0.165

contrary, ED and DP produced more visually smooth results due to their segment-based detection. Similar to Fig. 14b, we observed that, with comparable performance, ED is consistently faster over DP with different parameters.

Table 3 summarizes the comparison between ED and alternative methods in terms of running time, F1-Event and F1 scores averaged over sequences in the entire dataset. As what we have observed in Fig. 14, STBB had the smallest running time yet with the worst performance. Among the top performing DP and ED, without losing much accuracy, ED improved the speed against DP from about 6x to 14x.

7 Conclusion and Future Work

Using Branch-and-Bound (B&B), we introduced an unsupervised approach to common event discovery in segments

of variable length. We derived novel bounding functions so that the B&B framework guarantees a globally optimal solution in an empirically efficient manner. With slight modifications, the B&B framework can be readily applied to common event discovery, synchrony discovery, video search, and supervised event detection. The searching procedure can be extended to discovery among multiple time series, of multiple commonalities, and can be accelerated with warm start and parallelism. We evaluated the effectiveness of the B&B framework in motion capture of deliberate whole-body behavior and in video of spontaneous facial behavior in interviews, small groups of young adults, and parent-infant face-to-face interaction.

Future work includes promoting the scalability of the proposed algorithm. Given current pairwise design, the computational complexity grows quadratically with the number of input sequences. One direction is to pursue parallelism, *i.e.*, compute pairwise bounds independently using clusters or multi-threading, and then aggregate these bounds into a overall score.

Acknowledgements This work was supported in part by US National Institutes of Health grants GM105004 and MH096951. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Institutes of Health. The authors would like to thank Feng Zhou and Jiabei Zeng for helping partial experiments.

References

- Amberg, B., & Vetter, T. (2011). Optimal landmark detection using shape models and branch and bound. In *ICCV*.
- Balakrishnan, V., Boyd, S., & Balemí, S. (1991). Branch and bound algorithm for computing the minimum stability degree of parameter-dependent linear systems. *International Journal of Robust and Nonlinear Control*, 1(4), 295–317.
- Barbič, J., Safonova, A., Pan, J. Y., Faloutsos, C., Hodgins, J. K., & Pollard, N. S. (2004). Segmenting motion capture data into distinct behaviors. In *Proceedings of Graphics Interface 2004* (pp. 185–194). Canadian Human-Computer Communications Society.
- Bartlett, M. S., Littlewort, G. C., Frank, M. G., Lainscsek, C., Fasel, I. R., & Movellan, J. R. (2006). Automatic recognition of facial actions in spontaneous expressions. *Journal of Multimedia*, 1(6), 22–35.
- Begum, N., & Keogh, E. (2014). Rare time series motif discovery from unbounded streams. *VLDB*, 8(2), 149–160.
- Boiman, O., & Irani, M. (2005). Detecting irregularities in images and in video. In *ICCV*.
- Brand, M., Oliver, N., & Pentland, A. (1997). Coupled HMMs for complex action recognition. In *CVPR*.
- Brendel, W., & Todorovic, S. (2011). Learning spatiotemporal graphs of human activities. In *Proceedings of ICCV* (pp. 778–785).
- Chaaaraoui, A. A., Climent-Pérez, P., & Flórez-Revuelta, F. (2012). A review on vision techniques applied to human behaviour analysis for ambient-assisted living. *Expert Systems with Applications*, 39(12), 10873–10888.
- Chu, W. S., Chen, C. P., & Chen, C. S. (2010). Momi-cosegmentation: Simultaneous segmentation of multiple objects among multiple images. In *Proceedings of ACCV*.
- Chu, W. S., De la Torre, F., & Cohn, J. F. (2016). Selective transfer machine for personalized facial expression analysis. *TPAMI*.
- Chu, W. S., Zeng, J., De la Torre, F., Cohn, J. F., & Messinger, D. S. (2015). Unsupervised synchrony discovery in human interaction. In *ICCV*.
- Chu, W. S., Zhou, F., & De la Torre, F. (2012). Unsupervised temporal commonality discovery. In *ECCV*.
- Cooper, H., & Bowden, R. (2009). Learning signs from subtitles: A weakly supervised approach to sign language recognition. In *CVPR*.
- De la Torre, F., Chu, W. S., Xiong, X., Ding, X., & Cohn, J. F. (2015). Intraface. In *Automatic face and gesture recognition*.
- Delaherche, E., Chetouani, M., Mahdhaoui, A., Saint-Georges, C., Viaux, S., & Cohen, D. (2012). Interpersonal synchrony: A survey of evaluation methods across disciplines. *IEEE Transactions on Affective Computing*, 3(3), 349–365.
- Ding, X., Chu, W. S., De la Torre, F., Cohn, J. F., & Wang, Q. (2012). Facial action unit event detection by cascade of tasks. In *ICCV* (vol. 2013).
- Du, S., Tao, Y., & Martinez, A. M. (2014). Compound facial expressions of emotion. *Proceedings of the National Academy of Sciences*, 111(15), E1454–E1462.
- Duchenne, O., Laptev, I., Sivic, J., Bach, F., & Ponce, J. (2009). Automatic annotation of human actions in video. In *ICCV*.
- Everingham, M., Zisserman, A., Williams, C. I., & Van Gool, L. (2006). The PASCAL visual object classes challenge 2006 results. In *2th PASCAL challenge*.
- Feris, R., Bobbitt, R., Brown, L., & Pankanti, S. (2014). Attribute-based people search: Lessons learnt from a practical surveillance system. In *ICMR*.
- Gao, L., Song, J., Nie, F., Yan, Y., Sebe, N., & Tao Shen, H. (2015). Optimal graph learning with partial tags and multiple features for image and video annotation. In *CVPR*.
- Gendron, B., & Crainic, T. G. (1994). Parallel branch-and-branch algorithms: Survey and synthesis. *Operations Research*, 42(6), 1042–1066.
- Girard, J. M., Cohn, J. F., Jeni, L. A., Lucey, S., & De la Torre, F. (2015). How much training data for facial action unit detection? In *AFGR*.
- Goldberger, J., Gordon, S., & Greenspan, H. (2003). An efficient image similarity measure based on approximations of kl-divergence between two gaussian mixtures. In *ICCV*.
- Gusfield, D. (1997). *Algorithms on strings, trees, and sequences: Computer science and computational biology*. Cambridge: Cambridge University Press.
- Han, D., Bo, L., & Sminchisescu, C. (2009). Selection and context for action recognition. In *ICCV* (2009).
- Hoai, M., Zhong Lan, Z., & De la Torre, F. (2011). Joint segmentation and classification of human actions in video. In *CVPR*.
- Hongeng, S., & Nevatia, R. (2001). Multi-agent event recognition. In *ICCV*.
- Hu, W., Xie, N., Li, L., Zeng, X., & Maybank, S. (2011). A survey on visual content-based video indexing and retrieval. *IEEE Transactions on Systems, Man, and Cybernetics, Part C*, 41(6), 797–819.
- Jhuang, H., Serre, T., Wolf, L., & Poggio, T. (2007). A biologically inspired system for action recognition. In *ICCV*.
- Keogh, E., & Ratanamahatana, C. A. (2005). Exact indexing of dynamic time warping. *Knowledge and Information Systems*, 7(3), 358–386.
- Krüger, S. E., Schafföner, M., Katz, M., Andelic, E., & Wendenmuth, A. (2005). Speech recognition with support vector machines in a hybrid system. In *Interspeech*.
- Lampert, C., Blaschko, M., & Hofmann, T. (2009). Efficient subwindow search: A branch and bound framework for object localization. *IEEE TPAMI*, 31(12), 2129–2142.
- Laptev, I., Marszałek, M., Schmid, C., & Rozenfeld, B. (2008). Learning realistic human actions from movies. In *Proceedings of CVPR*.
- Lehmann, A., Leibe, B., & Van Gool, L. (2011). Fast prism: Branch and bound hough transform for object class detection. *IJCV*, 94(2), 175–197.
- Littlewort, G., Bartlett, M. S., Fasel, I., Susskind, J., & Movellan, J. (2006). Dynamics of facial expression extracted automatically from video. *Image and Vision Computing*, 24(6), 615–625.
- Liu, C. D., Chung, Y. N., & Chung, P. C. (2010). An interaction-embedded hmm framework for human behavior understanding: With nursing environments as examples. *IEEE Transactions on Information Technology in Biomedicine*, 14(5), 1236–1246.
- Liu, H., & Yan, S. (2010). Common visual pattern discovery via spatially coherent correspondences. In *Proceedings of CVPR*.
- Liu, J., Shah, M., Kuipers, B., & Savarese, S. (2011). Cross-view action recognition via view knowledge transfer. In: *CVPR*.
- Lucey, P., Cohn, J. F., Kanade, T., Saragih, J., Ambadar, Z., & Matthews, I. (2010). The extended cohn-kanade dataset (CK+): A complete dataset for action unit and emotion-specified expression. In: *CVPRW*.
- Maier, D. (1978). The complexity of some problems on subsequences and supersequences. *Journal of the ACM*, 25(2), 322–336.
- Matthews, I., & Baker, S. (2004). Active appearance models revisited. *IJCV*, 60(2), 135–164.
- Messinger, D. M., Ruvolo, P., Ekas, N. V., & Fogel, A. (2010). Applying machine learning to infant interaction: The development is in the details. *Neural Networks*, 23(8), 1004–1016.
- Messinger, D. S., Mahoor, M. H., Chow, S. M., & Cohn, J. F. (2009). Automated measurement of facial expression in infant-mother interaction: A pilot study. *Infancy*, 14(3), 285–305.
- Minnen, D., Isbell, C., Essa, I., & Starner, T. (2007). Discovering multivariate motifs using subsequence density estimation. In: *AAAI*.

- Mueen, A., & Keogh, E. (2010). Online discovery and maintenance of time series motifs. In: *KDD*.
- Mukherjee, L., Singh, V., & Peng, J. (2011). Scale invariant cosegmentation for image groups. In *Proceedings of CVPR*.
- Murphy, K. P. (2012). *Machine learning: A probabilistic perspective*. Cambridge: MIT press.
- Narendra, P. M., & Fukunaga, K. (1977). A branch and bound algorithm for feature subset selection. *IEEE Transactions on Computers*, 100(9), 917–922.
- Nayak, S., Duncan, K., Sarkar, S., & Loeding, B. (2012). Finding recurrent patterns from continuous sign language sentences for automated extraction of signs. *Journal of Machine Learning Research*, 13(1), 2589–2615.
- Oliver, N. M., Rosario, B., & Pentland, A. P. (2000). A bayesian computer vision system for modeling human interactions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8), 831–843.
- Paterson, M., & Dančák, V. (1994). Longest common subsequences. *Mathematical Foundations of Computer Science, 1994*(841), 127–142.
- Platt, J., et al. (1999). Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. *Advances in Large Margin Classifiers*, 10(3), 61–74.
- Reddy, K. K., & Shah, M. (2013). Recognizing 50 human action categories of web videos. *Machine Vision and Applications*, 24(5), 971–981.
- Rubner, Y., Tomasi, C., & Guibas, L. J. (2000). The earth mover's distance as a metric for image retrieval. *IJCV*, 40(2), 99–121.
- Sadanand, S., & Corso, J. J. (2012). Action bank: A high-level representation of activity in video. In *CVPR*.
- Sanginetto, E., Zen, G., Ricci, E. & Sebe, N. (2014). We are not all equal: Personalizing models for facial expression analysis with transductive parameter transfer. In *Proceedings of ACM MM*.
- Sayette, M. A., Creswell, K. G., Dimoff, J. D., Fairbairn, C. E., Cohn, J. F., Heckman, B. W., et al. (2012). Alcohol and group formation a multimodal investigation of the effects of alcohol on emotion and social bonding. *Psychological Science*, 23, 869–878.
- Schindler, G., Krishnamurthy, P., Lublinerman, R., Liu, Y., & Dellaert, F. (2008). Detecting and matching repeated patterns for automatic geo-tagging in urban environments. In *Proceedings of CVPR*.
- Schmidt, R. C., Morr, S., Fitzpatrick, P., & Richardson, M. J. (2012). Measuring the dynamics of interactional synchrony. *Journal of Nonverbal Behavior*, 36(4), 263–279.
- Scholkopf, B. (2001). The kernel trick for distances. In *NIPS*.
- Schuller, B., & Rigoll, G. (2006). Timing levels in segment-based speech emotion recognition. In *Interspeech*.
- Si, Z., Pei, M., Yao, B., & Zhu, S. (2011). Unsupervised learning of event and-or grammar and semantics from video. In *ICCV*.
- Sivic, J., & Zisserman, A. (2003). Video google: A text retrieval approach to object matching in videos. In *Proceedings of ICCV*.
- Sun, M., Telaprolu, M., Lee, H., & Savarese, S. (2012). An efficient branch-and-bound algorithm for optimal human pose estimation. In *CVPR*.
- Turaga, P., Veeraraghavan, A., & Chellappa, R. (2009). Unsupervised view and rate invariant clustering of video sequences. *CVIU*, 113(3), 353–371.
- Valstar, M., & Pantic, M. (2006). Fully automatic facial action unit detection and temporal analysis. In *CVPRW*.
- Viola, P., & Jones, M. J. (2004). Robust real-time face detection. *IJCV*, 57(2), 137–154.
- Wang, H., Zhao, G., & Yuan, J. (2014). Visual pattern discovery in image and video data: A brief survey. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 4(1), 24–37.
- Wang, Y., Jiang, H., Drew, M. S., Li, Z., & Mori, G. (2006). Unsupervised discovery of action classes. In *Proceedings of CVPR*.
- Wang, Y., & Velipasalar, S. (2009). Frame-level temporal calibration of unsynchronized cameras by using Longest Consecutive Common Subsequence. In *ICASSP*.
- Yang, Y., Saleemi, I., & Shah, M. (2013a). Discovering motion primitives for unsupervised grouping and one-shot learning of human actions, gestures, and expressions. *TPAMI*, 35(7), 1635–1648.
- Yang, Y., Song, J., Huang, Z., Ma, Z., Sebe, N., & Hauptmann, A. G. (2013b). Multi-feature fusion via hierarchical regression for multimedia analysis. *IEEE Transactions on Multimedia*, 15, 572–581.
- Yu, X., Zhang, S., Yu, Y., Dunbar, N., Jensen, M., Burgoon, J. K., & Metaxas, D. N. (2013). Automated analysis of interactional synchrony using robust facial tracking and expression recognition. In *Automatic Face and Gesture Recognition*.
- Yuan, J., Liu, Z., & Wu, Y. (2011). Discriminative video pattern search for efficient action detection. *IEEE TPAMI*, 33(9), 1728–1743.
- Zheng, Y., Gu, S., & Tomasi, C. (2011). Detecting motion synchrony by video tubes. In *ACMMM*.
- Zhou, F., De la Torre, F., & Hodgins, J. K. (2013). Hierarchical aligned cluster analysis for temporal clustering of human motion. *IEEE TPAMI*, 35(3), 582–596.
- Zhou, F., De la Torre, F., & Cohn, J. F. (2010). Unsupervised discovery of facial events. In *Proceedings of CVPR*.
- Zhu, S., & Mumford, D. (2006). A stochastic grammar of images. *Foundations and Trends in Computer Graphics and Vision*, 2(4), 259–362.