# Unsupervised Discriminative Embedding for Sub-Action Learning in Complex Activities

Sirnam Swetha[*], Hilde Kuehne[†], Yogesh S Rawat[*], Mubarak Shah[*]
[*]Center for Research in Computer Vision, University of Central Florida, Orlando, FL
[†]MIT-IBM Watson Lab, Cambridge, MA

## Abstract

*Action recognition and detection in the context of long untrimmed video sequences has seen an increased attention from the research community. However, annotation of complex activities is usually time consuming and challenging in practice. Therefore, recent works started to tackle the problem of unsupervised learning of sub-actions in complex activities. This paper proposes a novel approach for unsupervised sub-action learning in complex activities. The proposed method maps both visual and temporal representations to a latent space where the sub-actions are learnt discriminatively in an end-to-end fashion. To this end, we propose to learn sub-actions as latent concepts and a novel discriminative latent concept learning (DLCL) module aids in learning sub-actions. The proposed DLCL module lends on the idea of latent concepts to learn compact representations in the latent embedding space in an unsupervised way. The result is a set of latent vectors that can be interpreted as cluster centers in the embedding space. The latent space itself is formed by a joint visual and temporal embedding capturing the visual similarity and temporal ordering of the data. Our joint learning with discriminative latent concept module is novel which eliminates the need for explicit clustering. We validate our approach on three benchmark datasets and show that the proposed combination of visual-temporal embedding and discriminative latent concepts allow to learn robust action representations in an unsupervised setting.*

## 1. Introduction

Recent years have seen a great progress in video activity analysis. However, most of this research is focused on the classification of short video clips with atomic or short-range actions [1, 2, 3]. This is a relatively easier task when compared with analysis of untrimmed and complex video sequences [4, 5, 6, 7, 8, 9, 10, 11, 12]. In untrimmed video analysis, the focus is either on the problem of temporal *action localization* [13, 14, 15, 16], where only a set of key actions is considered in untrimmed videos; or on the task of
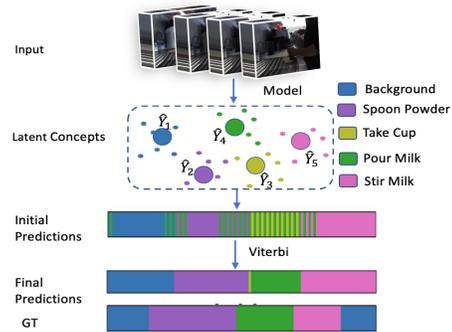


Figure 1: Overview of the proposed approach. Given videos of a complex activity, the proposed model learns sub-actions as latent concepts in an end-to-end manner. The latent concept assignment for each input video segment feature forms sub-action prediction shown as 'Initial Predictions', which is then refined using Viterbi to generate 'Final Predictions'. Sample results for activity 'Make Chocolate Milk', it can be seen that the latent concepts are able to group sub-actions. The sub-action 'pour-milk' includes lifting bottle and pouring milk; the jitter can be associated to the confusion when either a chocolate/milk bottle is lifted.

temporal *action segmentation* [9, 10, 17, 18, 19, 20], where each frame of the video is associated with a respective sub-action class as it requires to identify sub-actions and also temporally localize them.

Existing works on temporal action segmentation mainly explore supervised approaches where frame-level annotations are required for all the sub-actions [21, 8, 12, 9, 22, 11, 23]. However, complex activities are usually long-ranged and obtaining frame-level annotation is arduous and expensive. A new line of research focuses on learning these sub-actions from videos of a complex activity in an *unsupervised setting* [4, 5, 19, 10, 24]. In the unsupervised setting, the problem is even more challenging as it requires (i) breaking down a complex activity video into semantically meaningful sub-actions; and (ii) capturing the temporal relationship between the sub-actions. Most approaches tackle this problem in two stages, where during the first stage an

arXiv:2105.00067v1 [cs.CV] 30 Apr 2021

embedding based on visual and/or temporal information is learned, and in the second stage clustering is applied on this embedding space. This limits the learning ability by preventing the embedding to actually learn from clustering. At the same time, performing explicit clustering which is independent of embedding learning makes the model less efficient and does not utilize end-to-end learning.

To address this problem, we propose an end-to-end approach where sub-actions are learned by combining embedding and latent concepts. Here, the embedding space is trained jointly with the latent concepts leading to an effective sub-action discovery as shown in Figure 1. To allow for such a joint training, we propose a novel discriminative latent concept learning (DLCL) module which combines latent concept learning with a contrastive loss to ensure that the sub-actions learnt in the latent embedding are distant. The resulting latent concept vectors can be interpreted as cluster centers, removing the need for explicit clustering at a later stage.

As the sub-actions are softly bound to the temporal position of each activity, incorporating temporal ordering is crucial. Recent works [10, 19] incorporated temporal embedding either by predicting the discrete temporal entities or by learning continuous temporal embedding with shallow MLP architectures. In those cases, the temporal information is only given by a discrete or continuous scalar value and the joint embedding space is constructed by predicting this value from the input. To learn better spatio-temporal representations, we propose to use temporal position encoding [25] instead of scalar values and learn the respective embedding space by jointly reconstructing both visual and temporal representations. This embedding is further trained jointly with constrastive loss of the latent concept module, so that the embedding is also guided by and contributes to overall clustering.

We evaluate our method on three benchmark datasets: Breakfast [8], 50Salads [26] and YouTube Instructions [5]. For the evaluation at test time, we follow the protocol from [19] and employ the Viterbi algorithm to decode the initial sub-action predictions into coherent segments based on the ordered clustering of the sub-action latent concepts. A detailed analysis shows the impact of the proposed elements, the reconstruction and well as the latent concept learning.

In summary, we propose a novel end-to-end unsupervised approach for sub-action learning in complex activities. We make the following contributions in this work:

- We propose an unsupervised end-to-end approach for sub-action learning in complex activities by jointly learning an embedding which simultaneously incorporates visual and temporal information.

- We learn discriminative latent concepts using contrastive loss, thus integrating clustering as part of latent embedding learning.

- Our method improves the state-of-the-art on three benchmark datasets.

## 2. Related Work

Recently, there has been a lot of interest in learning with less supervision. This is essential for both action [3, 1, 2] and complex activity understanding [12, 27, 11], as supervised approaches require a large number of frames to be annotated in videos, which is expensive, tedious and cannot be scaled to large datasets. Weakly supervised approaches use a video and readily available information like accompanying text narration or audio. Some works [28, 29] use associated text narrations or scripts for learning actions in the video. Another line of work with weak-supervision include the works where it is assumed that the order of sub-actions is known [17, 30, 31, 32], however the per-frame annotations between video and sub-actions are not known during training. Authors in [33] propose to use combination of audio, text and video to identify steps in instructional videos in kitchen setting. The performance of the above methods is highly dependent on both the availability and quality of the text/audio alignment to video, which is not guaranteed and heavily limit their application.

There have been some works, in which the assumption of weak supervision have been removed in learning of action classes. One of the first works with no supervision addressed the problem of human motion segmentation [34] based on sensory-motor data, and proposed an application of a parallel synchronous grammar, to learn simple action representations similar to words in language. Later, a Bayesian non-parametric approach to concurrently model multiple sets of time series data was proposed in [35]. However, this work only focuses on motion capture data. [36, 37] benefit from the temporal structure of videos to fine-tune networks without any labels. Additionally, [38, 39, 40, 41] also leverage the temporal structure of videos to learn feature representation to learn actions.

Recently, unsupervised approaches have been proposed to learn sub-actions in complex activity. [10, 19, 24] propose unsupervised approaches for temporal segmentation of complex activity into sub-actions. While [4] proposes to solve a variant of the problem where the goal is to detect event boundaries, *i.e.* event boundary segmentation for complex activities. This does not focus on identifying sub-actions, instead it learns to identify boundaries across multiple sub-actions in long videos. A self-supervised predictive learning framework is proposed to solve by utilizing the difference between observed and predicted frame features to determine event boundaries in complex activities.

In this work, we focus on solving the temporal segmentation of complex activity into sub-actions as shown in [10, 19, 24]. In [10], an iterative approach is proposed that alternates between discriminative learning and gener-
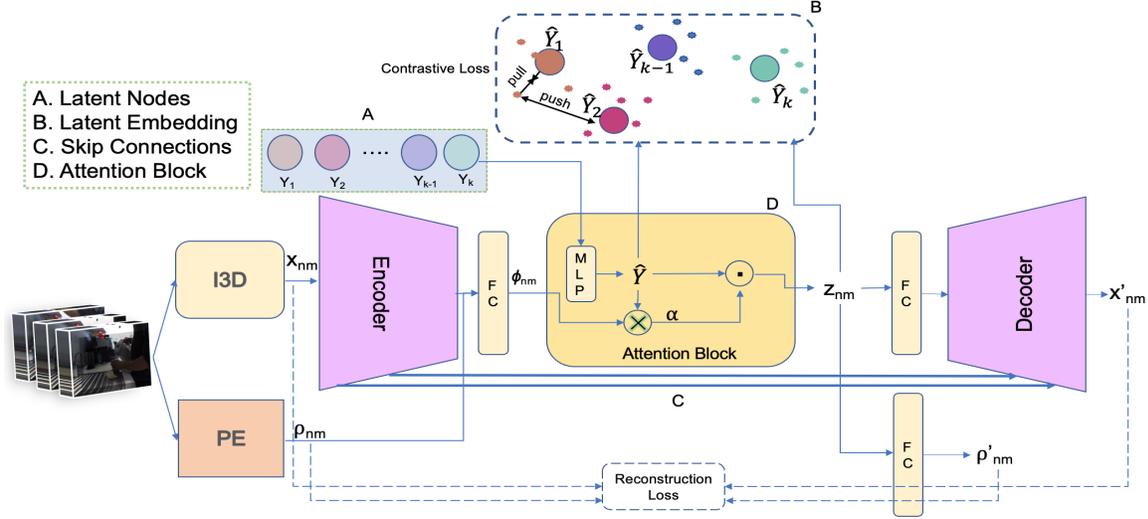
Figure 2: Overview of the proposed model. Given videos for a complex activity, we extract visual features $(x_{nm})$ and compute positional encoding vectors $(\rho_{nm})$ which are fed to the encoder to map them to a joint latent embedding for learning sub-action clusters. To learn these sub-action clusters as latent concepts $\widehat{Y}$, an attention block $(D)$ is used which takes in randomly initialized vectors $(Y)$ along with $\phi_{nm}$ and learns the latent concepts. We use contrastive loss to learn $\widehat{Y}$ discriminatively in $B$. Here $\alpha$, $z_{nm}$ and $\widehat{Y}_k$ represent attention, latent vector for input $x_{nm}$ and $k^{th}$ latent concept respectively.

ative modeling. For discriminative learning, they map the visual features into latent space using a 'linear' transformation and compute the weight matrix which minimizes the pairwise ranking loss. For temporal ordering they use Generalized Mallows Model which models distributions over permutations as they formulate complex activity as a sequence of permutable sub-actions. In [19], the model incorporates the continuous temporal ordering of frames in a joint embedding space. This is achieved by training a regressor to predict the timestamp of frames in a video. The hidden layer representations are used as the embedding for clustering and the clusters are ordered with respect to their time stamp. We refer to this model as CTE (Continuous Temporal Embedding). In [24], two-stage embedding pipeline is proposed where a next frame prediction U-Net model in stage one is combined with with temporal discriminator in stage two followed by clustering. The temporal embedding model employed is similar to [19].

Latent embedding learning is crucial for unsupervised learning, recently [7] formulated learning graph based latent embedding using latent concepts for supervised classification of complex activities. The intuition was to model long range videos using latent concepts as graphical nodes for complex activity recognition. Inspired by their ideology of latent concept learning to model latent space, we propose DLCL as an unsupervised latent learning module with joint embedding learning to model sub-actions.

Most of the above works in unsupervised learning involve two stage process which does not utilize end-to-end learning making them less efficient as clustering is inde-

pendent of embedding learning. In this work, we present an end-to-end model where clustering is incorporated in embedding learning using a constrastive loss. To incorporate temporal ordering we propose to use positional encodings and we also propose an effective way to unify visual and temporal representations to learn a visual-temporal embedding by jointly reconstructing visual and temporal representations. The proposed latent embedding is not only better at capturing visual & temporal representations but also clustering friendly. We demonstrate later in this paper the usefulness of the proposed model both qualitatively and quantitatively.

## 3. Proposed Model

### 3.1. Overview

Given a set of $N$ videos, $\{V_n\}_{n=1}^{N}$, for a complex activity, we divide each video into segments and for each segment we extract I3D features [1], and compute positional encoding vectors [25] as described in Section 3.2. Each video is represented by $M_n$ features where $x_{nm}$ represent the $m^{th}$ feature in the $n^{th}$ video, and its corresponding positional encoding is represented by $\rho_{nm}$. The task is to learn the sub-actions and their ordering for each activity, i.e., by predicting sequence of a sub-action labels $l_{nm} \in \{1, 2, ..., K\}$ for each feature $x_{nm}$ for each video. The number of sub-actions labels $K$ for each activity is the maximum number of possible sub-actions as they occur in that activity.

Overview of our proposed model is shown in Figure 2.

3

First we learn an encoded representation of $\mathrm{x}_{nm}$ and $\rho_{nm}$ shown as $\phi_{nm}$, which is passed as input feature to the 'Attention Block' (shown as $D$ in Figure 2) to learn the latent concepts/clusters which are representative of sub-actions. The attention block learns the latent concepts $\widehat{Y} \in \{\widehat{Y}_1, \widehat{Y}_2, ..., \widehat{Y}_K\}$ discriminatively, where each input feature ($\phi_{nm}$) is assigned to only one latent concept. We use a combination of reconstruction loss and constrastive loss to learn the embedding(shown as $B$ in Figure 2). We believe that using a combination of both visual and temporal information in conjunction with latent concept learning to learn a latent embedding (shown as block B in fig. 2) makes our model more robust. We evaluate the performance of our model based on latent concept assignments for each input feature which forms 'Initial Predictions'. Then, we model the sub-action transitions and perform Viterbi decoding to estimate optimal sub-action ordering.

Note that unlike previous works [19], we do not perform explicit clustering, instead our model learns to cluster features in latent space as discussed in Section 3.2 & 3.3. Thus eliminating the need for all the data to be available at once, our model can learn the latent concepts incrementally. The resulting sub-action latent concepts are temporally ordered and then each video is decoded w.r.t the above ordering given initial sub-action probability assignments for each clip to each latent concept as described in Section 3.5.

## 3.2. Joint Visual and Temporal Latent Embedding

In unsupervised learning, the approach to learn clusters in latent space plays a critical role in learning semantically meaningful clusters. We employ an encoder-decoder model to obtain the latent representation. Skip-connections are included between encoder and decoder (shown as $C$ in Figure 2), as they help to preserve commonality of an action and reduce redundant information like background in latent representation.

For incorporating temporal ordering in our model, we employ the positional encodings inspired by [25]. We divide the video segment sequence into g equal groups and then use the ordering index to compute positional encoding vectors. Quantizing temporal index of the video clip and using a positional encoding not only captures relative positioning but also makes it easy to generalize for highly varying video lengths. The idea of learning a mapping from features to joint visual and temporal embedding with an encoder-decoder aids in grouping clips into sub-actions in the latent space. The reconstruction loss for the auto-encoder is composed of visual features and positional encoding as shown below,

$$Loss_r = L(\mathrm{x}_{nm}, \mathrm{x}'_{nm}) + \beta * L(\rho_{nm}, \rho'_{nm}), \quad (1)$$

where, $\mathrm{x}_{nm}$, $\mathrm{x}'_{nm}$ respectively represent input and reconstructed visual feature; $\rho_{nm}$, $\rho'_{nm}$ respectively represent in-

put and reconstructed positional encoding; $\beta$ is a hyperparameter and $L$ is a loss function penalizing $\mathrm{x}'_{nm}$ and $\rho'_{nm}$ for being dissimilar from $\mathrm{x}_{nm}$ and $\rho_{nm}$ respectively, namely mean squared error. A combination of latent visual feature representation and the positional encodings becomes input to the 'Attention Block'. In order to ensure that the learnt clusters are representative of sub-actions, the clusters have to be distant in the latent space, which is described in the next section.

## 3.3. Discriminative Latent Concept Learning

The idea behind having this module is to learn the sub-action clusters discriminatively in the latent space in an end-to-end fashion, eliminating the need for explicit clustering. The attention block is inspired by [7], which takes an input feature ($\phi_{nm}$) and randomly initialized latent vectors ($Y$) which is analogous to cluster center initialization as shown in Figure 2. The latent concepts ($\widehat{Y}$) are learnt using an MLP with weight ($w$) and bias ($b$) i.e., it transforms the random latent vector initializations ($Y$) to latent concepts ($\widehat{Y}$) as $\widehat{Y} = w * Y + b$. Though latent vector initialization ($Y$) is fixed, $w$ & $b$ are learnable parameters making the latent concepts ($\widehat{Y}$) learnable in the latent space. These latent concepts which represent cluster centers are learned by minimizing the contrastive loss by moving features in the latent space closer to the latent concepts. The similarity between input feature ($\phi_{nm}$) and latent concepts ($\widehat{Y}$) is measured with the dot product $\otimes$. Then, activation function $\sigma$ is applied on the similarities to compute activation values $\alpha$ i.e., $\alpha = \sigma(\phi_{nm} * \widehat{Y}^T)$. Finally, the attended latent vector representation is computed as $Z_{nm} = \alpha \odot \widehat{Y}$, which captures how much each latent concept is related to the given input feature. However, these latent concepts tend to learn similar/overlapping concepts, which is not what we intend to learn. Our objective in learning the latent concepts is to cluster the latent representations discriminatively. We achieve this with a contrastive loss, where the similarity between latent vectors of the same sub-action with the maximum confident latent concept is maximized, while the similarity w.r.t other latent concepts is minimized as shown in Eq 2.

$$Loss_d(Z_{nm}, \widehat{Y}) = -log \frac{e^{sim(Z_{nm}, \widehat{Y}_{k^*})}}{\sum_{k \neq k^*} e^{sim(Z_{nm}, \widehat{Y}_k)}} \quad (2)$$

where, $\widehat{Y}_k$ represents the latent concept associated with $k^{th}$ sub-action, $sim$ denotes cosine similarity and $k^*$ represents the latent concept with maximum confidence probability for $Z_{nm}$ as shown below

$$k^* = \underset{k}{\mathrm{argmax}} P(k|Z_{nm}) \quad (3)$$

where $P(k|Z_{nm}) = \sigma(sim(Z_{nm}, \widehat{Y}_k))$ represents the confidence probability of latent vector $Z_{nm}$ for the latent concept $\widehat{Y}_k$, $\sigma$ is softmax activation.
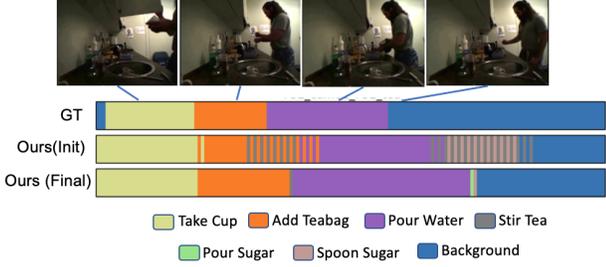
4

Figure 3: Qualitative comparison of initial predictions (w/o viterbi) and after Viterbi predictions of our approach for activity 'Make Tea'. It can be seen that our model ('Init') is able to group sub-actions and also learn the ordering of sub-actions for an activity. The jitter in sub-action predictions occurs during transition from one sub-action to next, which is expected during transition. Finally, using transition modeling Viterbi decoding smoothness the jitter between sub-action transitions.

## 3.4. Overall Loss.

Total loss for learning the proposed embedding is composed of losses from Section 3.2 and 3.3 as $Loss = \lambda * Loss_r + \gamma * Loss_d$

## 3.5. Temporal Segmentation

**Initial Predictions** At test time, we first assign each feature in video to its respective closest latent concept vector using Eq 3. This gives initial predictions directly based on the embedding (shown as predictions in Figure 1). For ease of understanding, we refer to those as latent sets, analogous to clusters, from here on.

**Transition modeling and Viterbi decoding**
Figure 4 represents a brief outline for transition modeling and Viterbi decoding. To allow for a temporal decoding, the global ordering of the latent sets needs to be estimated. We follow here the protocol proposed by [19] and compute the mean timestamp for each set (shown as T in Figure 4) and sort them in ascending order. The last set in the sorted ordering becomes the terminal state and using this ordering the sub-action state transition probabilities from sub-action $i$ to $j$ are defined as $P(j|i)$ given:

$$P(j|i) = \begin{cases} 0.5, & \text{if } j \text{ immediately follows } i \\ 0.5, & \text{if } i = j \\ 1.0, & \text{if } i = j \ \& \ j \text{ is terminal state} \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

**Decoding** Finally, we use the ordering and transition probabilities to compute the best path for the set ordering given the input features $x_{nm}$ and $\rho_{nm}$. Using Eq. 3 we compute the probability of each embedded input feature ($Z_{nm}$) belonging to the latent set $k$. We maximize the probability of the input sequence following the order defined by Eq. 4
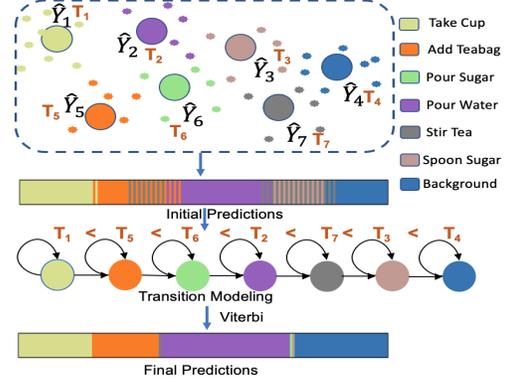


Figure 4: Brief overview of transition modeling and Viterbi decoding. Each latent concept is color coded (best viewed in color). The latent concepts are ordered w.r.t the mean time (shown as $T$) and each video is decoded into coherent segments using Viterbi algorithm based on the ordered sub-action latent concepts.

to get consistent latent set assignments in a video by maximizing,

$$\bar{l}_1^{M_n} = \operatorname*{argmax}_{l_1,...,l_{M_n}} \prod_{m=1}^{M_n} P(l_m|l_{m-1}) * P(l_m|Z_{nm}), \quad (5)$$

where $l_1,...,l_m \in \{1,2,...,K\}$ represent the set label sequence for $n^{th}$ video, $P(l_m = k|Z_{nm})$ is the probability that $Z_{nm}$ belongs to the $k^{th}$ latent set (as described in Section 3.3), $\bar{l}_1^{M_n}$ is the set label sequence for the maximum likelihood for $n^{th}$ video.
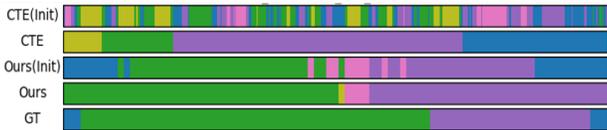
## 4. Experiments

For our experiments, we define a segment as a sequence of 8 frames. The video segment sequence is divided into 128 equal groups and then the ordering index is used to compute positional encoding [25] for each segment. We extract I3D features (layer 'mixed_5c') which is fed to the encoder. Our embedding dimension is 1024. We use a 3-layer encoder-decoder with Adam optimizer and the learning rate is set to $1 \times 10^{-4}$. We evaluate our approach on 3 datasets.
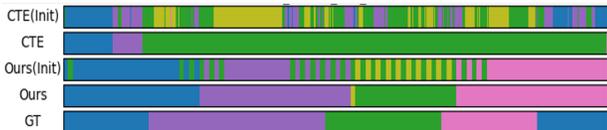**Breakfast Dataset** comprises of 10 complex activities of humans performing common kitchen activities. There are a total of 48 sub-activities in 1,712 videos with varying lengths based on activity and preparation style with variations in sub-action orderings.
**50Salads Dataset** contains videos of duration 4.5 hours for a single complex activity 'making mixed salad'. It is a multimodal dataset, as it includes RGB frames, depth maps and accelerometer data. However, we only use RGB frames. The videos in this dataset are much longer with average frame length of 10k frames and provides annotations at multiple granularity levels.
**YouTube Instructions Dataset** has 5 activities and 150 videos with 47 sub-actions. These videos are taken from

(a) Make Cereals



(b) Make Chocolate Milk

Figure 5: Illustrative comparison with state-of-the-art. By comparing with CTE (Init) and Ours (Init), we show that our approach learns to model sub-actions with very few intermittent sub-action transitions leading to effective grouping of sub-actions. Then, Viterbi decoding helps to smoothen the intermittent jitters in predictions. We show that our method provides coherent sub-action predictions and is able to capture the orderings for sub-actions.

| Method | F1-score | MoF |
|---|---|---|
| *Weakly Supervised* | | |
| RNN-FC [42] | | 33.3% |
| TCFPN [17] | | 38.4% |
| NN-Vit [32] | | 43% |
| D3TW [43] | | 45.7% |
| CDFL [44] | | 50.2% |
| *Unsupervised* | | |
| Mallow [10] | - | 34.6% |
| CTE [19] | 26.4% | 41.8% |
| JVT [24] | 29.9% | 48.1% |
| **Ours** | **31.9%** | 47.4% |
| LSTM-AL [4]* | - | 42.9%* |
| **Ours*** | - | **74.6%*** |

Table 1: Comparison of the proposed method to state-of-the-art on Breakfast dataset. Here, * denotes results with video-based Hungarian matching for the task event boundary segmentation.

YouTube directly and have background segments where there is no sub-action. The frequency and spread of background varies based on activity as well as on the person performing the task. Hence, the background segments neither have similar appearance nor have a temporal ordering. Therefore the background segments would be assigned to the latent concepts with very less confidence probability. Following protocol in [19], we consider $\tau$ percent of clips with least confidence as background. Only the foreground labeled segments along with latent concepts assignments form our initial predictions. We report results for background ratio of 60%.

| Method | MoC | MoC |
|---|---|---|
| | *w/o Viterbi* | *w Viterbi* |
| CTE [19] with FV | 20.9% | 40.1% |
| CTE [19] with I3D | 24.8% | 36.8% |
| **Ours with I3D** | **37.5%** | **46.9%** |

Table 2: Comparison of MoC (Mean over class) of all activities on Breakfast dataset before and after applying Viterbi. FV represents Fisher Vectors.

**Metrics** Our model predicts a sequence of cluster labels $\in \{1, 2, ..., K\}$ for each video without any correspondence to the $K$ ground-truth class labels. To map ground-truth and prediction label correspondences, inline with [5, 10, 19], for each activity we use the Hungarian method to find a one-to-one mapping for each cluster to exactly one sub-action and report performance after this mapping. In this work, we use Mean over Frames (MoF) as used by [10, 19] as well as F1-score used by [5]. In addition, we report Mean over class (MoC) accuracy, as it averages the accuracy for each activity class, therefore giving equal weights to all classes irrespective of the underlying data distribution. MoF is the percentage of correct predictions computed across all activity classes together, which can be affected by the underlying activity classes distribution and biased towards dominant activity class. For F1-score, similar to previous methods, we report the mean score over all activities. For state-of-the-art comparisons, we also evaluate our method for the task of event boundary segmentation following the protocol in [4] and compare our method to [4] - indicated as video-based Hungarian matching.

### 4.1. Comparison to state-of-the-art

Here, we compare the proposed method to state-of-the-art approaches. We present the accuracy comparison with recent works on Breakfast dataset in Table 1 and present the performance on new metric MoC in Table 2. Our approach achieves $47.4\%$ MoF and $31.9\%$ F1-score which is $2\%$ gain over state-of-the-art as shown in table 1. We show qualitative evaluation of the proposed approach in Figure 3 & 5. In Figure 5, we show that our approach models the sub-actions coherently with very less intermittent sub-action transitions along with learning ordering of sub-actions for complex activity. For example, in Figure 3 our model predicts 'stir-tea' with intermittent transitions after 'pour-water', this occurs when the person dips the tea bag in water which closely resembles to the sub-action 'stir-tea' (as shown in last image in Figure 3) and then it correctly predicts background once the dip action ends (there is no annotation for 'dipping tea-bag' in ground truth) indicating the goodness of the proposed sub-action learning. The intermittent transitions indicate that the model confuses to assign latent concept

| Method | F1-score | MoF |
|--------|----------|-----|
| CTE [19] | - | 35.5% |
| JVT [24] | - | 30.6% |
| **Ours** | **34.4%** | **42.2%** |
| LSTM-AL [4]* | - | 60.6*% |
| **Ours*** | - | **70.2%*** |

Table 3: Comparison of the proposed method to state-of-the-art unsupervised approaches on 50Salads dataset at granularity 'eval'. Here, * denotes results with video-based Hungarian matching for the task event boundary segmentation.

| Method | F1-score | MoF |
|--------|----------|-----|
| Frank-Wolfe [5] | 24.4% | 34.6% |
| Mallow [10] | 27.0% | 27.8% |
| CTE [19] | 28.3% | 39.0% |
| JVT [24] | 29.9% | 28.2% |
| **Ours** | 29.6% | **43.8%** |
| LSTM-AL [4]* | 39.7%* | - |
| **Ours*** | **45.4%*** | - |

Table 4: Comparison of the proposed method to state-of-the-art unsupervised methods on YouTube Instructions dataset. Here, * denotes results with video-based Hungarian matching for the task event boundary segmentation.

based on single feature and Viterbi aids in generating more coherent sub-action segments for the sequence as shown. Additionally, we evaluate our method for the task of event boundary segmentation and compare with the state-of-the-art approaches. Our approach out-performs the state-of-the-art MoF by a margin of 31% on Breakfast dataset indicating the effectiveness of the proposed method to temporally segment meaningful sub-actions.

For 50Salads dataset, we perform evaluation on granularity level 'eval' and provide state-of-the-art comparison in Table 3. Our method out-performs [19] by 6.67% and [24] by 11.6% with an F1-score of 34.37%. We further evaluate our method for the task of event boundary segmentation and perform state-of-the-art comparison in Table 3. We show 10% gain over state-of-the-art [4] MoF, indicating our method is effective in sub-action learning for complex events.

For YouTube Instructions dataset, we follow protocol in [5, 10, 19] and report the performance of our approach without considering the background frames. We achieve 42% MoC & 43.8% MoF (as shown in Table 4). This is a 4.8% gain in MoF over state-of-the-art method with comparable F1-score. Note that [4] reported F1-score with background frames included on YouTube Instructions Dataset. We follow the same procedure and compare our method to [4] in Table 4 (indicated with *). It can be seen that our method outperforms the state-of-the-art for event boundary segmentation task showing the sub-action learning capability to identify better event boundaries.

### 4.2. Evaluation of the Embedding.

To demonstrate the impact of the proposed embedding, we compare our Joint Embedding with Continuous Temporal Embedding in [19] in Table 2. From Table 2 (*MoC w/o viterbi*), it can be seen that the proposed joint embedding outperforms the continuous temporal embedding by a huge margin of 16.6%. It can be seen that our '*MoC w/o viterbi*' is closer to the CTE '*MoC w Viterbi*' suggesting that our embedding is very effective. To emphasize that our gain in

performance is due to the effectiveness of the approach and not with using I3D features, we train [19] using I3D features by keeping the embedding dimension same as ours and compare the performance. As shown in Table 2, the MoC w/o Viterbi improves by 4% by using I3D features on CTE, while the MoC with Viterbi drops by 3% with 1% increase in F1-score. However, our approach still outperforms the baseline (with same embedding dimension) by huge margin indicating our approach effectiveness.

Besides dataset level comparisons, we also show activity level comparison with CTE [19]. Figure 6 (a) shows that our joint embedding outperforms CTE on all activities indicating the significance of our joint embedding. We see a drop in performance for activity 'making cereals' after Viterbi decoding (from Figure 6(b)), this can be attributed to the ordering of the sub-actions 'take-bowl' and 'pour-cereals'. For many samples in 'making cereals', the sub-action 'take-bowl' does not occur impacting the ordering of both sub-actions leading to drop in performance.
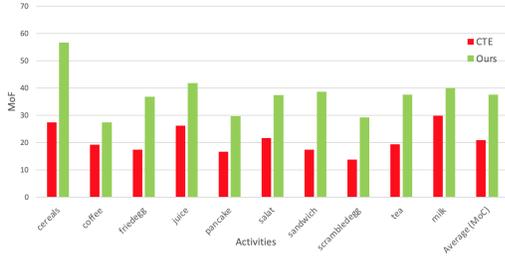
### 4.3. Ablation Experiments

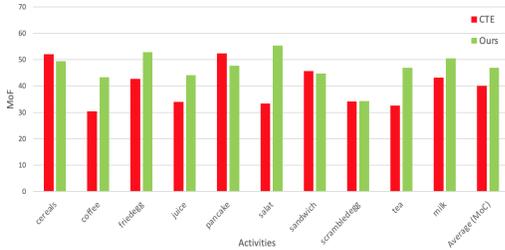We perform the below ablation studies on the breakfast dataset.

**Effect of Loss Components.** To begin with, we first examine the influence of $Loss_r$ and $Loss_d$ on our model and the performances are presented in Table 5. It can be seen that having all loss components leads to best performance.

**Effect of Discriminative learning.** The use of constrastive loss ($Loss_d$) helps the clusters to move apart in the latent space. This helps in obtaining more discrete boundaries in the latent space. As shown in Table 5, the accuracy drastically reduces to 35.8% (11% ↓) indicating the importance of discriminative learning.

**Effect of Positional Encoding.** Positional Encoding plays a crucial role in our model. It helps to temporally group the video clips in the latent space. As sub-actions are softly bound to the temporal position for each activity, removing

(a) MoF *w/o Viterbi*



(b) Final MoF

Figure 6: Activity level MoF comparison on Breakfast dataset with CTE [19]. Last column represents the average (MoC) for all activities. (a) represents MoF for each activity without Viterbi i.e, the MoF is computed based on the learnt cluster assignments. Our method outperforms the baseline on all activities. (b) represents MoF for each activity after applying Viterbi.

| $Loss_r$ | | $Loss_d$ | MoC |
|---|---|---|---|
| $L_f$ | $L_p$ | | |
| ✓ | - | - | 25.7% |
| - | ✓ | - | 33.6% |
| ✓ | ✓ | - | 35.8% |
| ✓ | - | ✓ | 40.2% |
| - | ✓ | ✓ | 40.1% |
| ✓ | ✓ | ✓ | 46.9% |

Table 5: Ablation experiments for the loss components are performed on the Breakfast dataset. $Loss_r$ and $Loss_d$ represents reconstruction loss and contrastive loss respectively. $L_f$ and $L_p$ denote the reconstruction loss for feature and positional encoding respectively.

reconstruction loss for positional encoding is expected to deteriorate the model performance. We observe the similar trend in Table 5. Additionally, we perform an ablation by removing the PE component branch and train our model end-to-end. As expected, there is a significant reduction in accuracy and F1-score (as shown in Table 6) indicating the significance of using positional encoding.

**Effect of Skip-Connections.** To assess the effectiveness of skip-connections, we report performance by removing the skip-connections and train model end-to-end. We re-
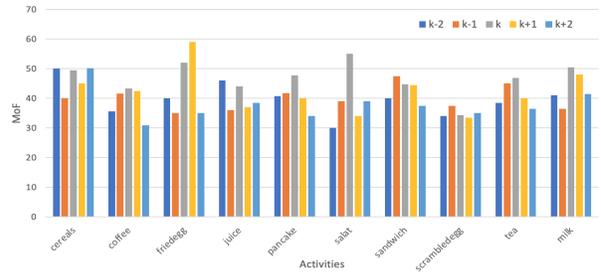


Figure 7: MoF vs. #sub-actions for all activities in Breakfast dataset. k represents the number of sub-actions from ground-truth; we vary the sub-actions for each activity and report MoF.

port the performance in Table 6, it can be seen that w/o skip-connections, the accuracy drops considerably indicating that the skip-connections help in learning better representations.

| | w/o PE | w/o SC | full |
|---|---|---|---|
| **MoC** | 40.9% | 35.7% | **46.9%** |
| **F1-score** | 20.3% | 28.7% | **31.9%** |

Table 6: Ablations experiments to evaluate the effect of PE and SC on Breakfast dataset (w/o: without, PE: Positional Encoding, SC: skip-connections).

**Effect of Sub-actions Cluster Size.** For all the above evaluations, the sub-action cluster size (K) is defined as mentioned in Section 3.1. To analyze the impact of sub-action cluster size, we vary the number of sub-actions from $K-2$ to $K+2$ where $K$ is the number of sub-actions as per ground truth and evaluate performance. Figure 7 shows the MoF vs number of sub-actions for each activity in Breakfast dataset. For 6 out of 10 activities we see that having K sub-actions leads to best performance.

## 5. Conclusion

In this work we proposed an end-to-end approach for unsupervised learning of sub-actions in complex activities. The main motivation behind this approach is to design a latent space to incorporate visual as well as positional encoding together. This latent space is learned via jointly training this embedding space in conjunction with a contrastive learning for clustering. We show that this allows for a robust learning that on it's own already results in a reasonable clustering of sub-actions. We then predict optimal sub-action sequence by employing the Viterbi algorithm which outperforms all the other methods. Our evaluation shows the impact of the proposed ideas and how they are able to improve the performance on this task compared to existing methods.

# References

[1] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *CVPR*, 2017.

[2] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *ICCV*, 2019.

[3] Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition in videos. In *NeurIPS*, 2014.

[4] Sathyanarayanan N Aakur and Sudeep Sarkar. A perceptual prediction framework for self supervised event segmentation. In *CVPR*, 2019.

[5] Jean-Baptiste Alayrac, Piotr Bojanowski, Nishant Agrawal, Ivan Laptev, Josef Sivic, and Simon Lacoste-Julien. Unsupervised learning from narrated instruction videos. In *CVPR*, 2016.

[6] Noureldien Hussein, Efstratios Gavves, and Arnold WM Smeulders. Timeception for complex action recognition. In *CVPR*, 2019.

[7] Noureldien Hussein, Efstratios Gavves, and Arnold WM Smeulders. Videograph: Recognizing minutes-long human activities in videos. *arXiv*, 2019.

[8] H. Kuehne, A. B. Arslan, and T. Serre. The language of actions: Recovering the syntax and semantics of goal-directed human activities. In *CVPR*, 2014.

[9] Hilde Kuehne, Juergen Gall, and Thomas Serre. An end-to-end generative framework for video segmentation and recognition. In *WACV*, 2016.

[10] Fadime Sener and Angela Yao. Unsupervised learning and segmentation of complex activities from video. In *CVPR*, 2018.

[11] Zheng Shou, Jonathan Chan, Alireza Zareian, Kazuyuki Miyazawa, and Shih-Fu Chang. Cdc: Convolutional-de-convolutional networks for precise temporal action localization in untrimmed videos. In *CVPR*, 2017.

[12] Serena Yeung, Olga Russakovsky, Greg Mori, and Li Fei-Fei. End-to-end learning of action detection from frame glimpses in videos. In *CVPR*, 2016.

[13] Peihao Chen, Chuang Gan, Guangyao Shen, Wenbing Huang, Runhao Zeng, and Mingkui Tan. Relation attention for temporal action localization. *IEEE Transactions on Multimedia*, 2019.

[14] Fuchen Long, Ting Yao, Zhaofan Qiu, Xinmei Tian, Jiebo Luo, and Tao Mei. Gaussian temporal awareness networks for action localization. In *CVPR*, 2019.

[15] Mengmeng Xu, Chen Zhao, David S Rojas, Ali Thabet, and Bernard Ghanem. G-tad: Sub-graph localization for temporal action detection. *arXiv*, 2019.

[16] Runhao Zeng, Wenbing Huang, Mingkui Tan, Yu Rong, Peilin Zhao, Junzhou Huang, and Chuang Gan. Graph convolutional networks for temporal action localization. In *ICCV*, 2019.

[17] Li Ding and Chenliang Xu. Weakly-supervised action segmentation with iterative soft boundary assignment. In *CVPR*, 2018.

[18] Alexander Richard, Hilde Kuehne, and Juergen Gall. Action sets: Weakly supervised action segmentation without ordering constraints. In *CVPR*, 2018.

[19] Anna Kukleva, Hilde Kuehne, Fadime Sener, and Jurgen Gall. Unsupervised learning of action classes with continuous temporal embedding. In *CVPR*, 2019.

[20] Yazan Abu Farha and Jurgen Gall. Ms-tcn: Multi-stage temporal convolutional network for action segmentation. In *CVPR*, 2019.

[21] Yazan Abu Farha and Jurgen Gall. Ms-tcn: Multi-stage temporal convolutional network for action segmentation. In *CVPR*, 2019.

[22] Hilde Kuehne, Alexander Richard, and Juergen Gall. A hybrid rnn-hmm approach for weakly supervised temporal action segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 2018.

[23] Colin Lea, Michael D. Flynn, Rene Vidal, Austin Reiter, and Gregory D. Hager. Temporal convolutional networks for action segmentation and detection. In *CVPR*, 2017.

[24] Rosaura G VidalMata, Walter J Scheirer, Anna Kukleva, David Cox, and Hilde Kuehne. Joint visual-temporal embedding for unsupervised learning of actions in untrimmed sequences. In *WACV*, 2020.

[25] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.

[26] Sebastian Stein and Stephen J McKenna. Combining embedded accelerometers with computer vision for recognizing food preparation activities. In *UBICOMP*, 2013.

[27] De-An Huang, Li Fei-Fei, and Juan Carlos Niebles. Connectionist temporal modeling for weakly supervised action labeling. In *ECCV*, 2016.

[28] Ivan Laptev, Marcin Marszalek, Cordelia Schmid, and Benjamin Rozenfeld. Learning realistic human actions from movies. In *CVPR*, 2008.

[29] Ozan Sener, Amir R Zamir, Silvio Savarese, and Ashutosh Saxena. Unsupervised semantic parsing of video collections. In *ICCV*, 2015.

[30] Hilde Kuehne, Alexander Richard, and Juergen Gall. Weakly supervised learning of actions from transcripts. *CVIU*, 2017.

[31] Alexander Richard, Hilde Kuehne, and Juergen Gall. Weakly supervised action learning with rnn based fine-to-coarse modeling. In *CVPR*, 2017.

[32] Alexander Richard, Hilde Kuehne, Ahsan Iqbal, and Juergen Gall. Neuralnetwork-viterbi: A framework for weakly supervised video learning. In *CVPR*, 2018.

[33] Jonathan Malmaud, Jonathan Huang, Vivek Rathod, Nick Johnston, Andrew Rabinovich, and Kevin Murphy. What's cookin'? interpreting cooking videos using text, speech and vision. *arXiv*, 2015.

[34] Gutemberg Guerra-Filho and Yiannis Aloimonos. A language for human action. *Computer*, 2007.

[35] Emily B Fox, Michael C Hughes, Erik B Sudderth, Michael I Jordan, et al. Joint modeling of multiple time series via the beta process with application to motion capture segmentation. *The Annals of Applied Statistics*, 2014.

[36] Xiaolong Wang and Abhinav Gupta. Unsupervised learning of visual representations using videos. In *ICCV*, 2015.

[37] Biagio Brattoli, Uta Buchler, Anna-Sophia Wahl, Martin E Schwab, and Bjorn Ommer. Lstm self-supervision for detailed behavior analysis. In *CVPR*, 2017.

[38] Vignesh Ramanathan, Kevin Tang, Greg Mori, and Li Fei-Fei. Learning temporal embeddings for complex video analysis. In *ICCV*, 2015.

[39] Basura Fernando, Efstratios Gavves, Jose M Oramas, Amir Ghodrati, and Tinne Tuytelaars. Modeling video evolution for action recognition. In *CVPR*, 2015.

[40] Anoop Cherian, Basura Fernando, Mehrtash Harandi, and Stephen Gould. Generalized rank pooling for activity recognition. In *CVPR*, 2017.

[41] Hsin-Ying Lee, Jia-Bin Huang, Maneesh Singh, and Ming-Hsuan Yang. Unsupervised representation learning by sorting sequences. In *ICCV*, 2017.

[42] Alexander Richard, Hilde Kuehne, and Juergen Gall. Weakly supervised action learning with rnn based fine-to-coarse modeling. In *CVPR*, 2017.

[43] Chien-Yi Chang, De-An Huang, Yanan Sui, Li Fei-Fei, and Juan Carlos Niebles. D3tw: Discriminative differentiable dynamic time warping for weakly supervised action alignment and segmentation. In *CVPR*, 2019.

[44] Jun Li, Peng Lei, and Sinisa Todorovic. Weakly supervised energy-based learning for action segmentation. In *ICCV*, 2019.