

# Recognize Human Activities from Partially Observed Videos

---

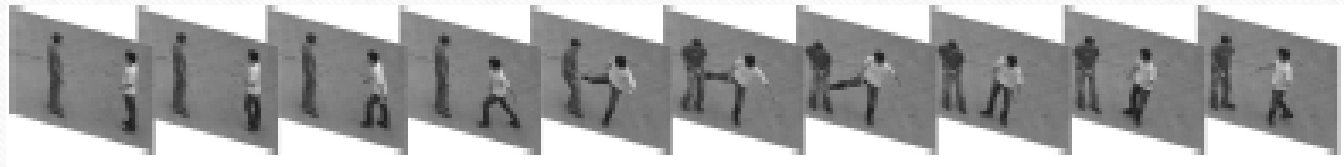
Anuroop Kakkirala   Praneeth AS



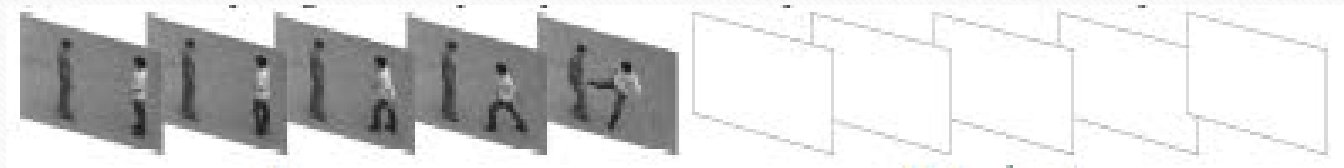
# Aim of the Project

---

- To recognize human activities from partially observed videos and predicting the unobserved subsequence in the video.
- Two possibilities:
  - 1) An Unobserved subsequence is at the end of the video.
  - 2) It may occur at any time by yielding a temporal gap in the video



Full Observation



Missing Observation  
At the end.



Gap filling



# Problem Formation:

## 1) For a Fully Observed Video

---

- Given a fully observed video  $O[1 : T]$  of length  $T$ ,  $O[t]$  indicates the frame at time  $t$ .
- Goal is to classify the video  $O[1 : T]$  into one of  $P$  activity classes  $\{A = \{A_p\}, p = 1, \dots, P\}$ .
- Human actions - sequence of simple actions - contain different spatiotemporal features. Divide uniformly into  $M$  different segments  $O(t_{i-1} : t_i]$  where  $t_i = i \cdot T/M$ ,  $i$  th stage of activity  $i = 1, \dots, M$ .

# Formation continued..

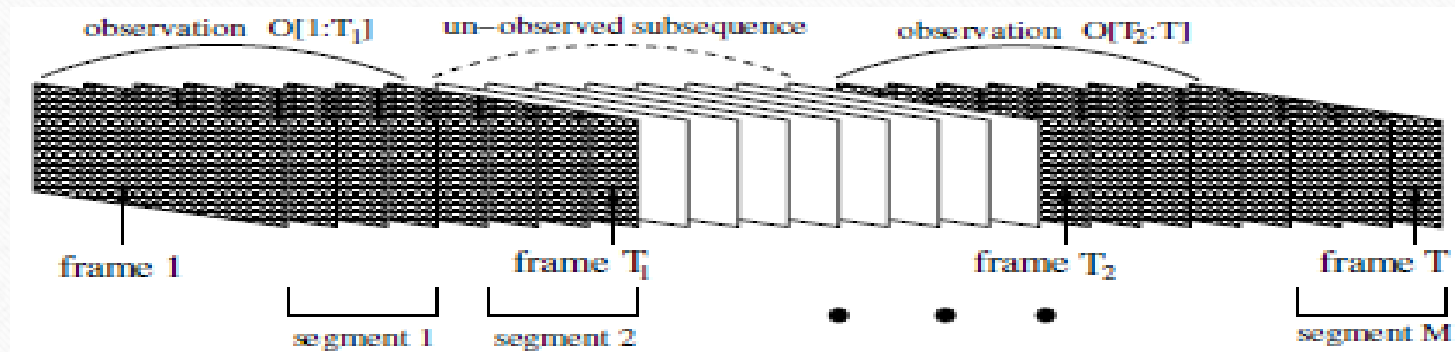
---

- $P(A_p | O[1 : T]) \propto \sum_{i|0 \text{ to } M} P(A_{p,(t_i-1,t_i)} | O[1 : T])$   
 $\propto \sum_{i|0 \text{ to } M} P(A_{p,(t_i-1,t_i)})P(O[1 : T] | A_{p,(t_i-1,t_i)}).$
- $P(A_{p,(t_i-1,t_i)}) = \text{prior of stage } i \text{ of activity } A_p$
- $P(O[1 : T] | A_{p,(t_i-1,t_i)}) = \text{observation likelihood given activity class } A_p.$
- $p^* = \arg \max \sum_{i|0 \text{ to } M} P(A_{p,(t_i-1,t_i)})P(O[1 : T] | A_{p,(t_i-1,t_i)})$
- $p^*$  is index of recognized activity.



## 2) For a Partially Observed Video

- Partially observed video -  $O[1 : T_1] \cup [T_2 : T]$ , where frames  $O(T_1 : T_2)$  are missing.
- $T_1$  is always the last frame of a segment and  $T_2$  is always the first of another segment.



## Formation continued..

---

- Posterior probability that an activity is presented in this partially observed video
- $P(A_p | O[1 : T_1] \cup [T_2 : T]) \propto w_1 \sum_{i | t_i \leq T_1} P(A_p, (t_i - 1, t_i] | O[1 : T_1]) + w_2 \sum_{i | t_i - 1 \geq T_2} P(A_p, (t_i - 1, t_i] | O[T_2 : T])$
- $w_1 = T_1 / (T_1 + T - T_2 + 1)$  ,  $w_2 = T - T_2 + 1 / (T_1 + T - T_2 + 1)$
- $p^* = \arg \max P(A_p | O[1 : T_1] \cup O[T_2 : T])$
- Where  $P(A_p | O[1 : T_1] \cup O[T_2 : T])$  can be calculated as above.



# Likelihood calculation

---

- Compare  $O[1 : T1]$  with the  $i$  th segment of all the training videos.
- Each segment of a video, use the bag-of-visual-words technique to organize its spatiotemporal features into a fixed-dimensional feature vector.
- $\mathbf{h}_i^n$  - feature(row) vector after applying bag-of-visual-words techniques to  $i$ th segment of the  $n$ th training video
- $\mathbf{h}_i^O$  - feature for stage  $i$  in  $O[1 : T1]$



# Likelihood calculation continued..

---

$$\bar{\mathbf{h}}_i = \frac{1}{N} \sum_{n=1}^N \mathbf{h}_i^n$$

$$P(\mathcal{O}[1 : T] | \mathcal{A}_p, (t_{i-1}, t_i]) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{-||h_i^{\mathcal{O}} - \bar{\mathbf{h}}_i||^2}{2\sigma^2}}$$