

---

# Aligning Videos In Space and Time

Senthil Purushwalkam<sup>1</sup>, Tian Ye<sup>1</sup>, Saurabh Gupta<sup>3</sup>, Abhinav Gupta<sup>1,2</sup>

<sup>1</sup>Carnegie Mellon University

<sup>2</sup>Facebook AI Research

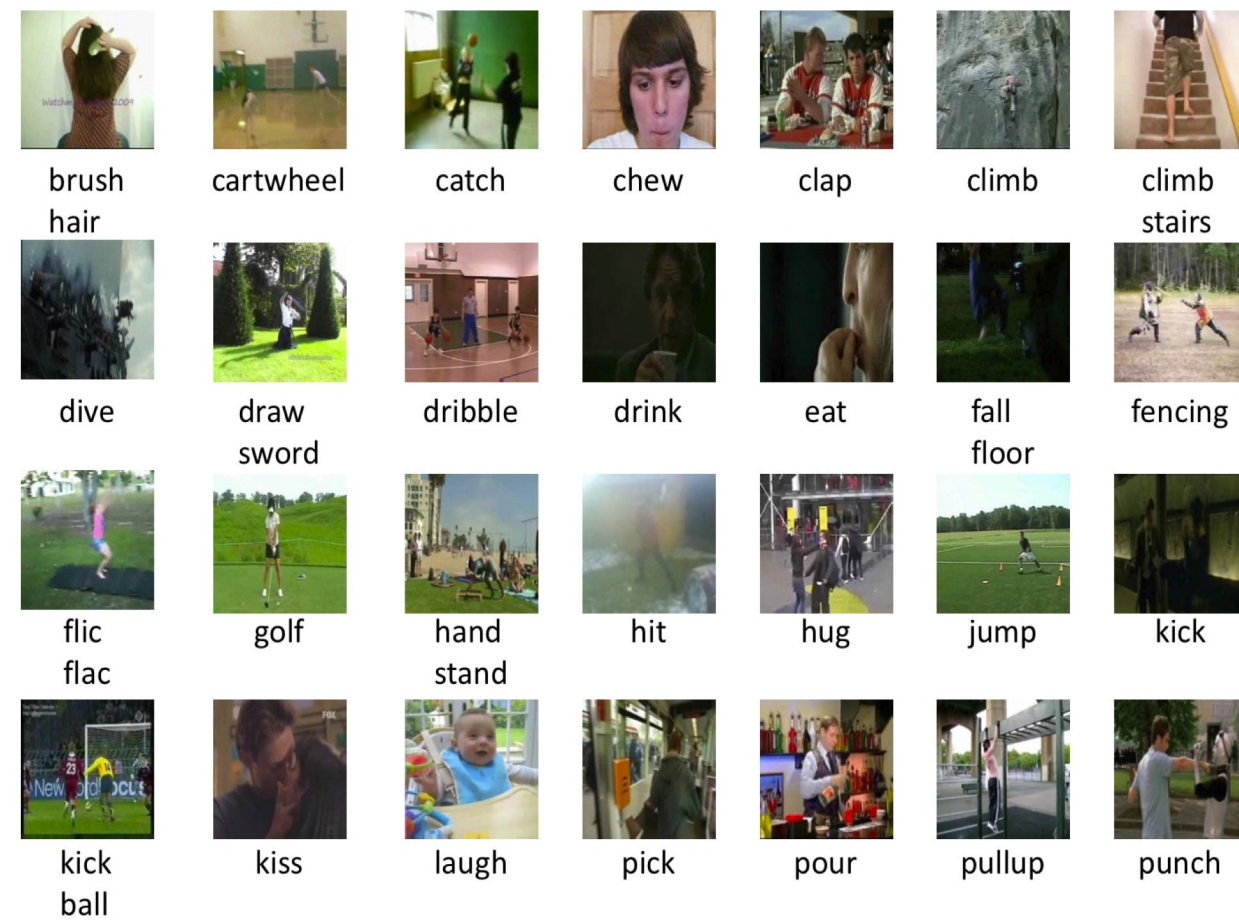
<sup>3</sup>University of Illinois Urbana-Champaign (UIUC)

European Conference on Computer Vision (ECCV), 2020

Short Presentation

---

# Video Understanding in Computer Vision



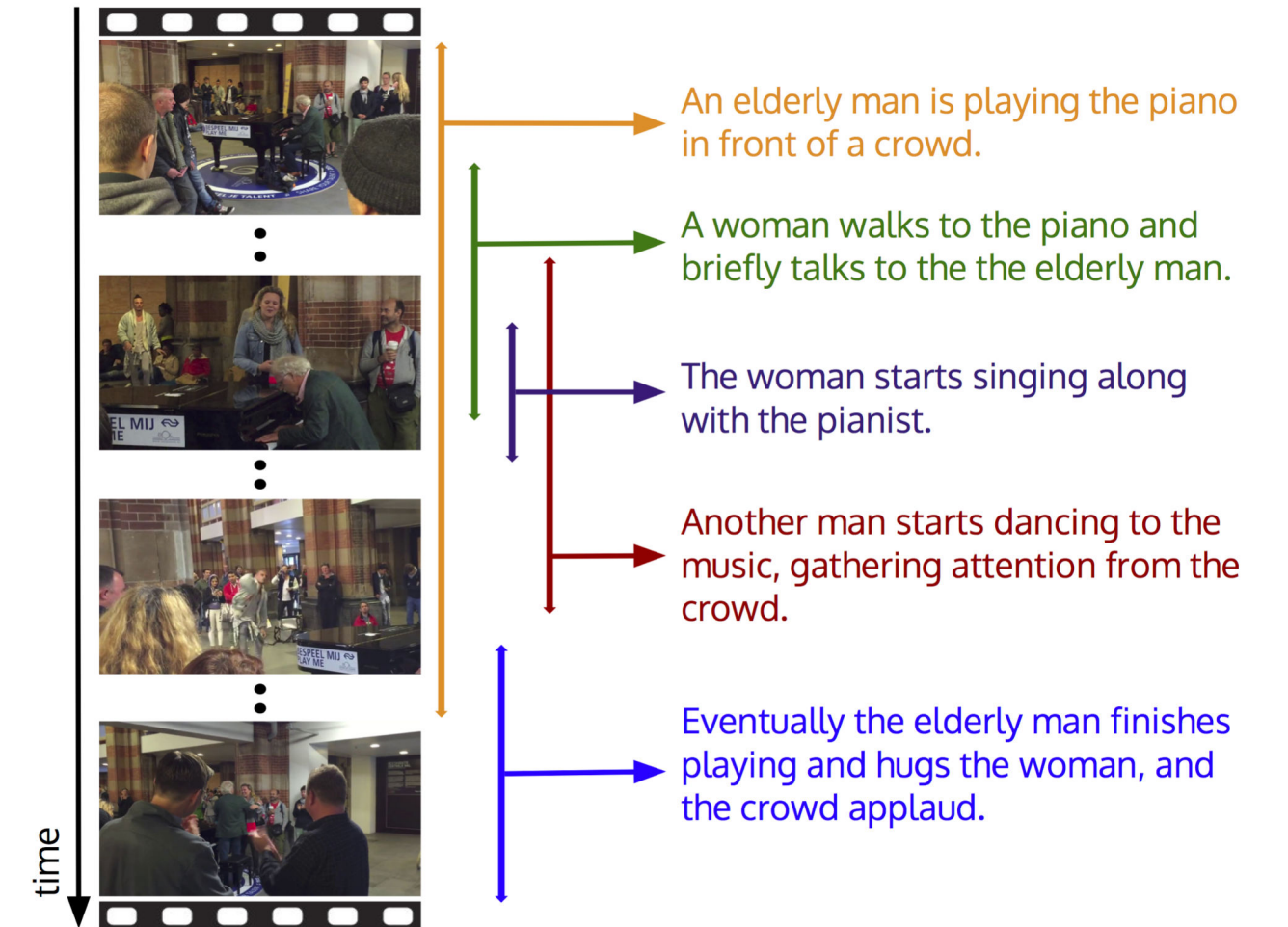
**Action Recognition**

Classification of Videos into  
Predefined Action Categories



**Action Detection**

Localizing Predefined Actions  
Temporally in Videos



**Video Captioning**

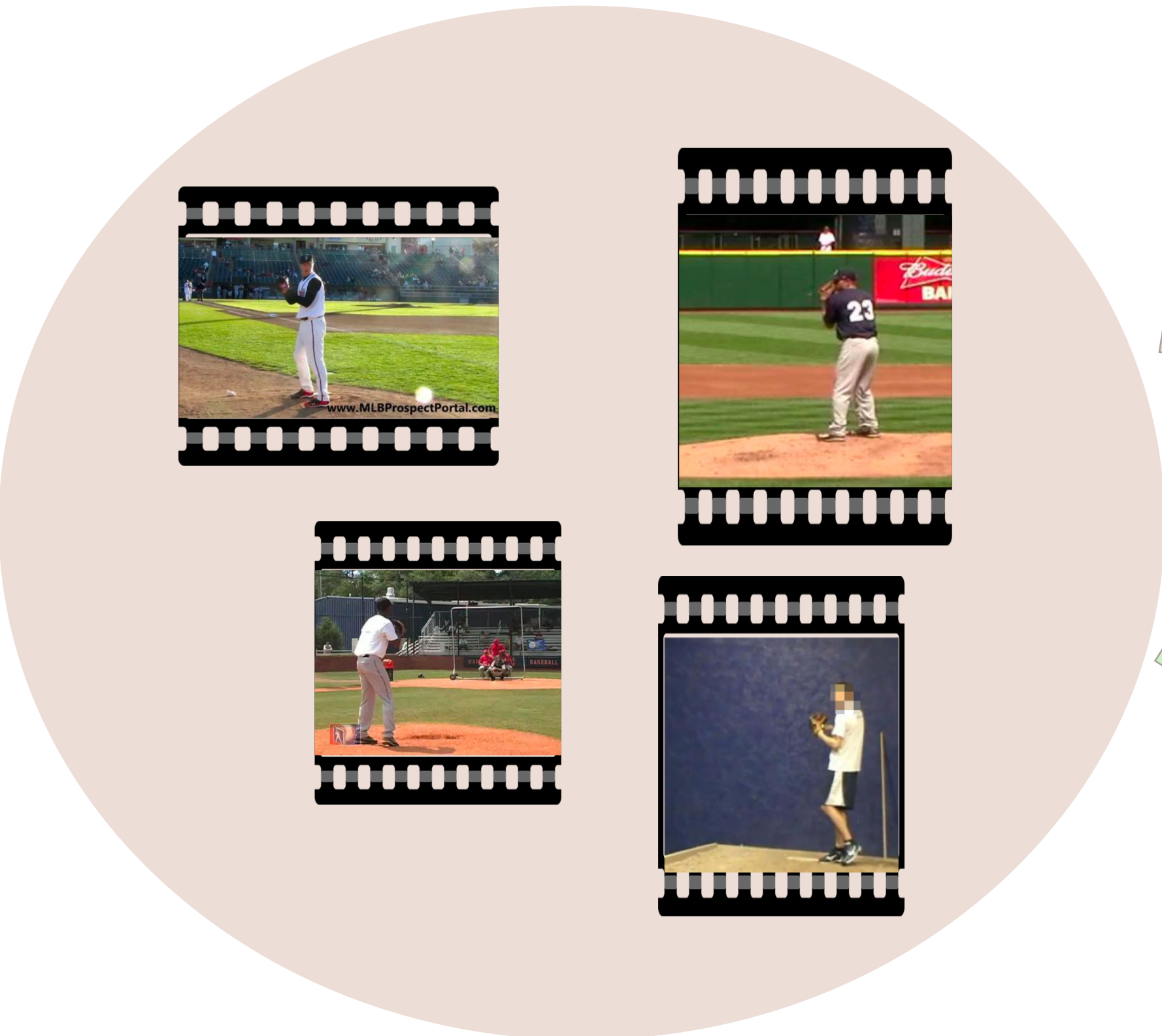
Generating Textual  
Descriptions for Videos

- Coarse Understanding of Videos
- Data Collection is Not Scalable to Denser Annotations



# Video Understanding via Association

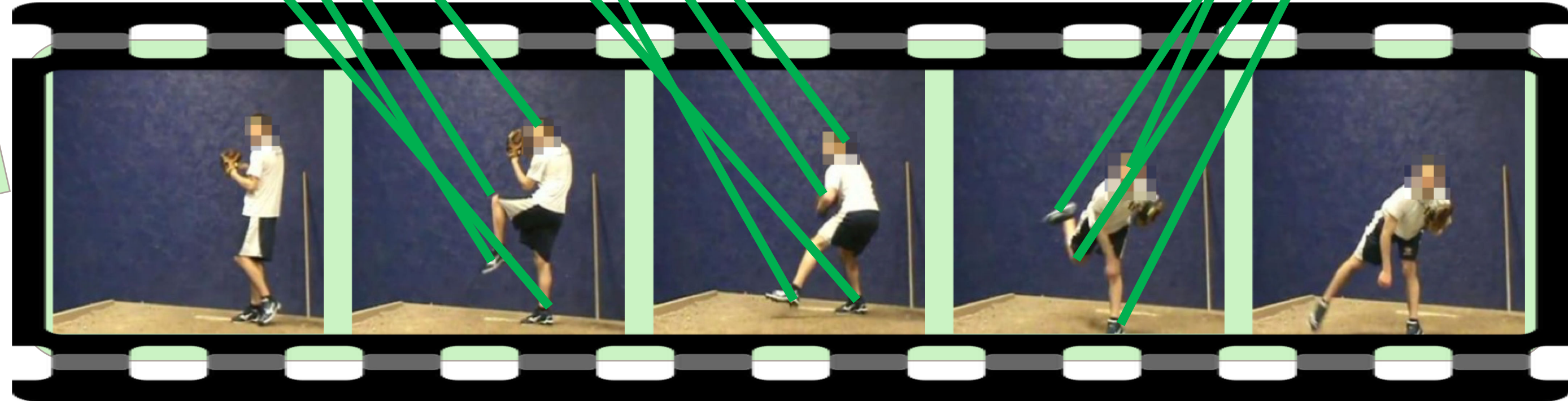
*Ask not "what is this?", ask "what is this like".  
-Moshe Bar*



Baseball Bowling



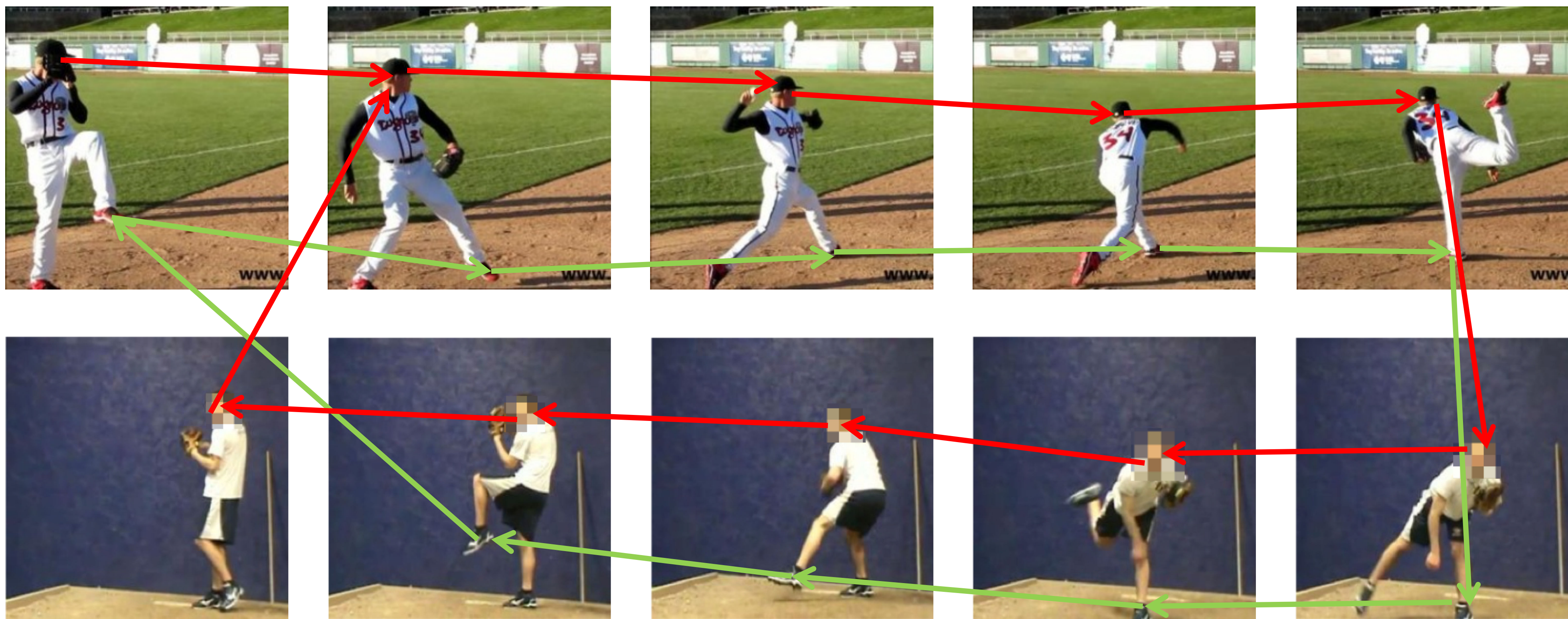
Query Video



Retrieved Video



# Unsupervised Learning of Spatio-Temporal Associations



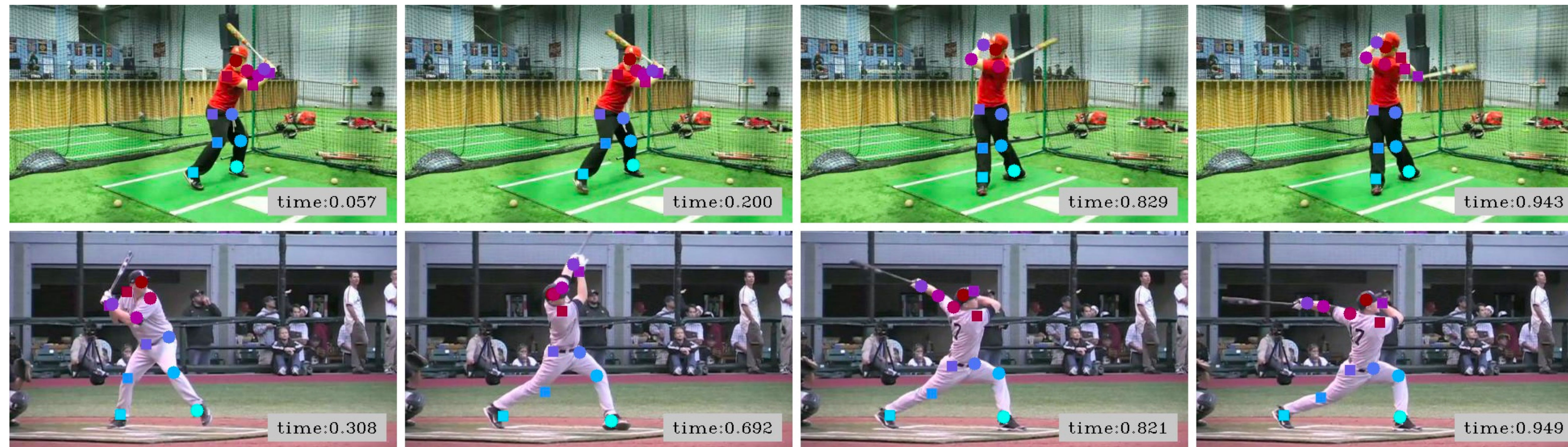
Positive Cycle

Negative Cycle



# Qualitative Results

Learned representation can effectively *spatio-temporally align videos*



Nearest neighbor using learned representation shows that it *encodes object states and appearance*



Thank You for Listening!

Project Webpage: <http://www.cs.cmu.edu/~spurushw/publication/alignvideos/>