

Self-supervised learning (for videos)

Computer Vision (SJK02)

Universitat Jaume I

The issues with labelling

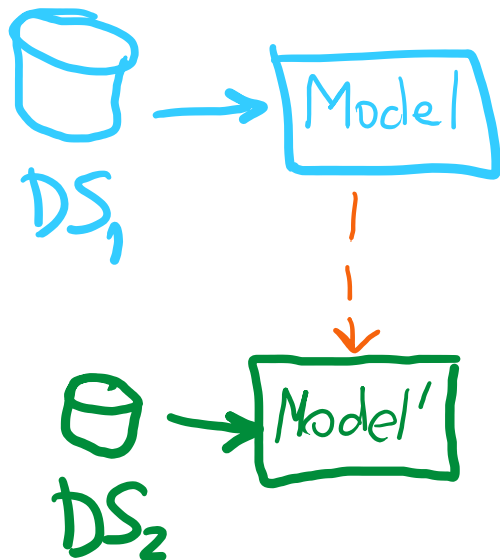
Annotation cost

Annotation bias

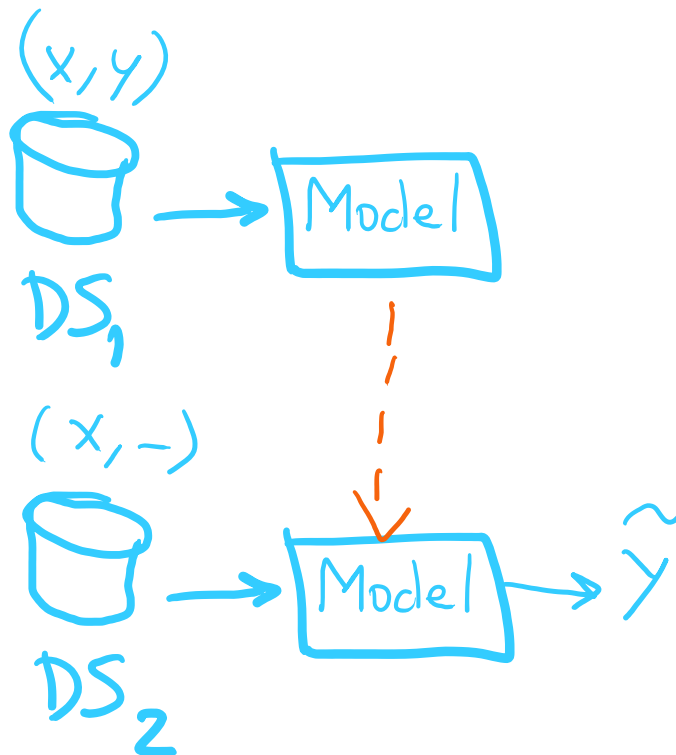
(Lack of domain generalisation)

(Lack of robustness)

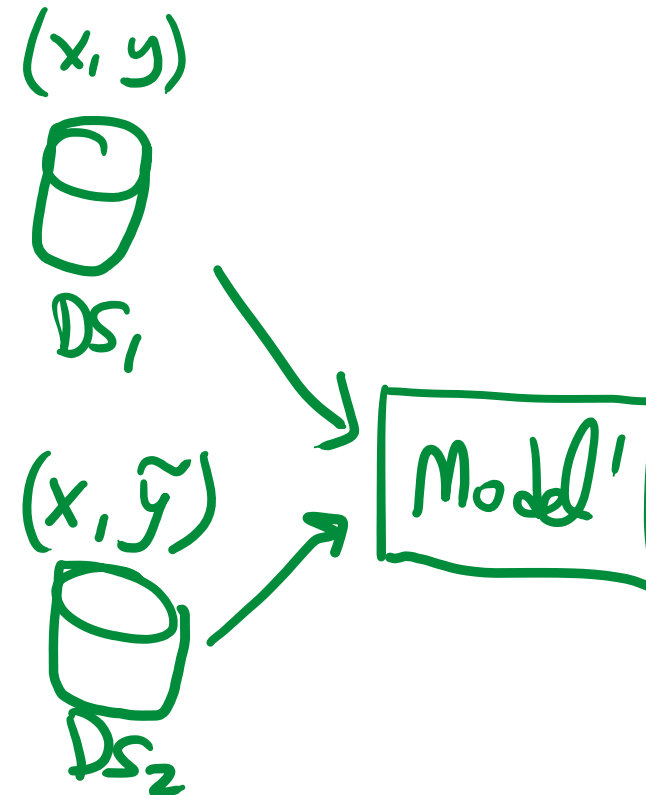
Pre-training \neq self-training \neq self-supervision



Pre-training
Transfer learning
Fine-tuning



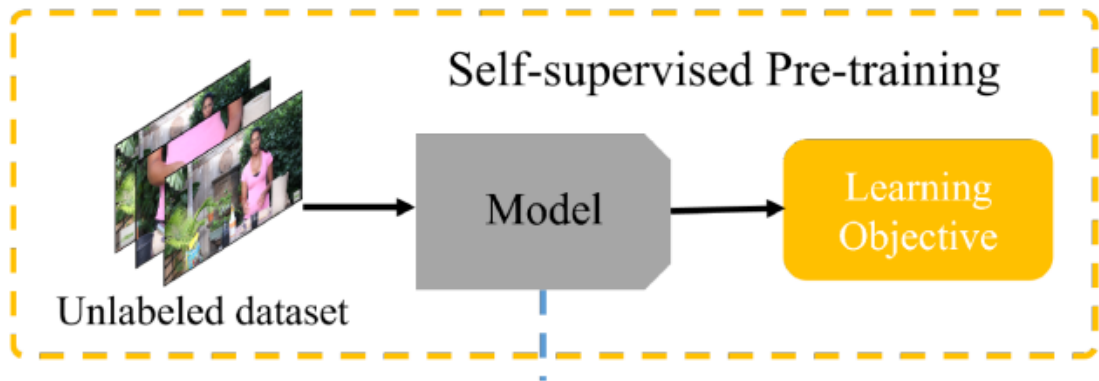
Self-training
Labeled + Pseudo-labeled datasets



Self-supervised learning (SSL) in a nutshell

Alternative to pretraining a model

- Large dataset (without labels!)
- Results in higher generalisation



Downstream tasks

Action **recognition**

Temporal action **segmentation**

Temporal Action **Step Localization**

Video **retrieval**

Text-to-Video Retrieval

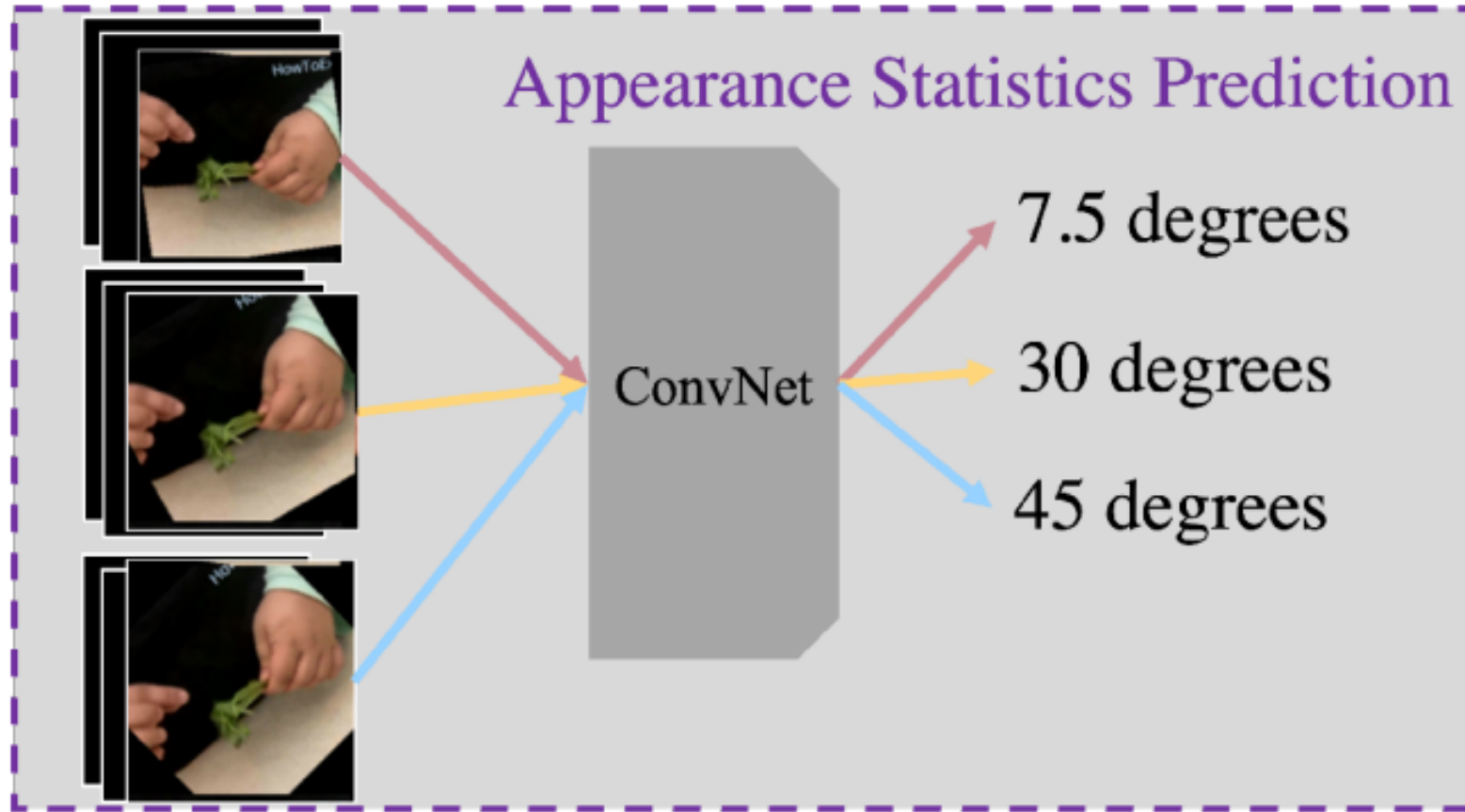
Video **Captioning**

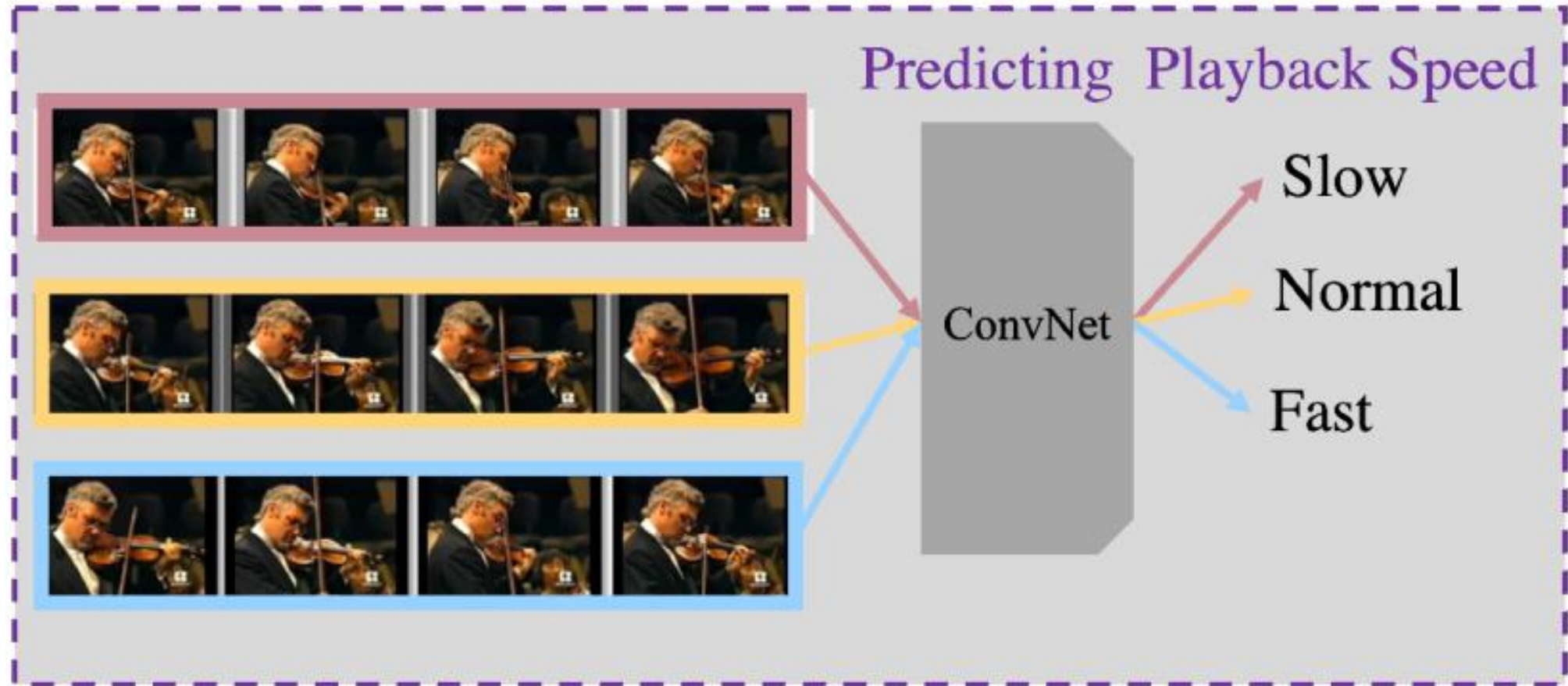
Approaches

Pretext
Generative
Contrastive
Multimodal*

(*) Not seen here

Pretext tasks

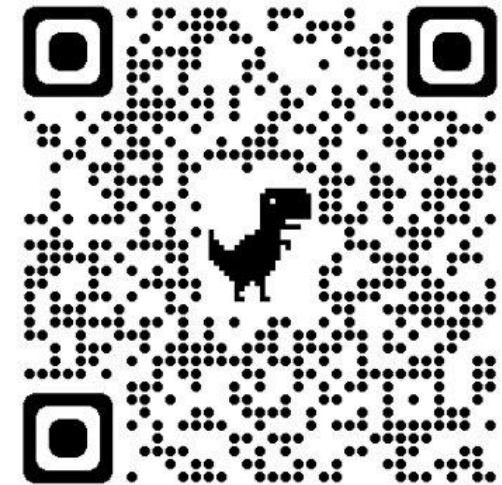




Why is this approach more specific to videos?

Other pretext task: *your turn*

Think of other pretext task for SSL for videos



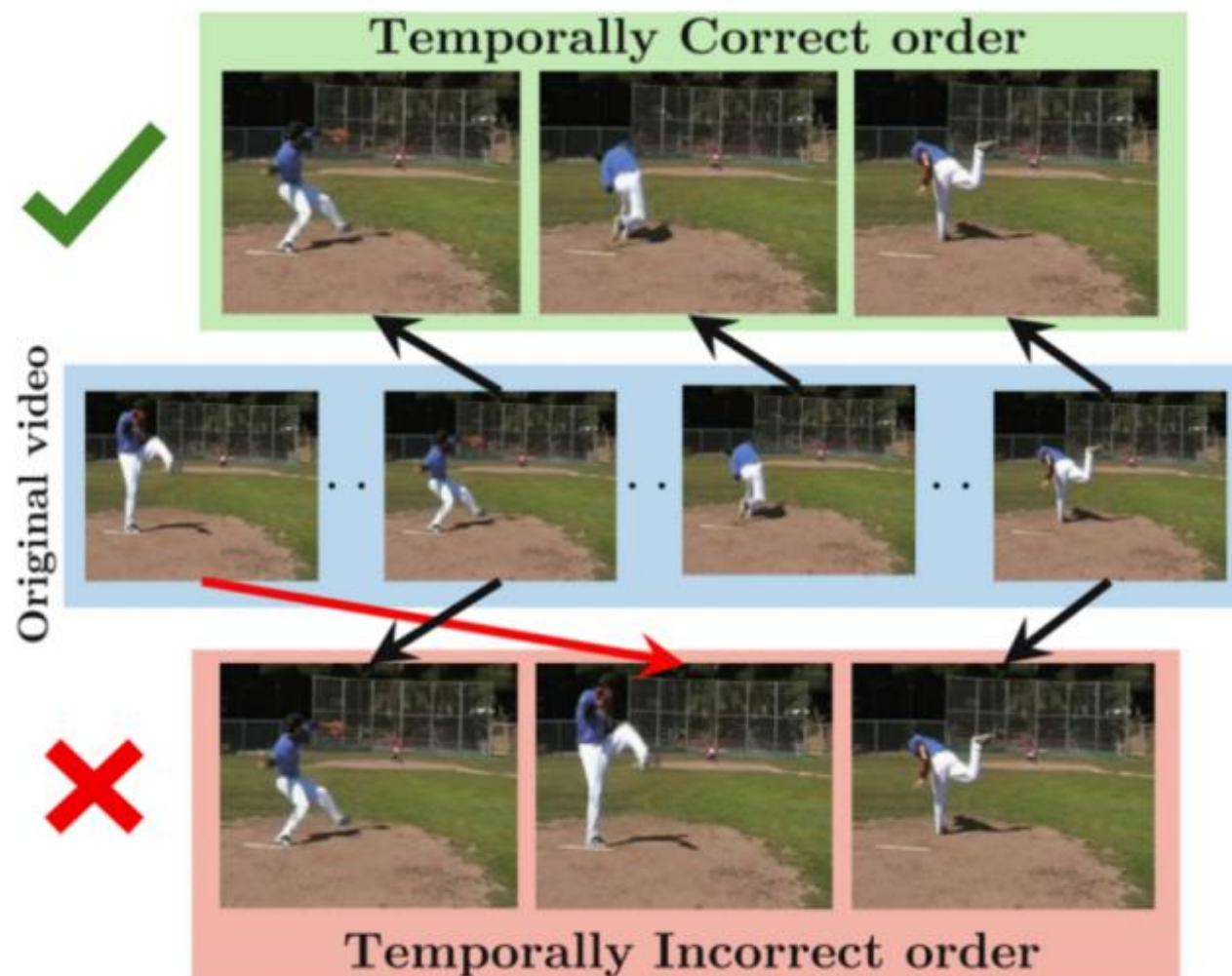
www.socrative.com

Room 219986

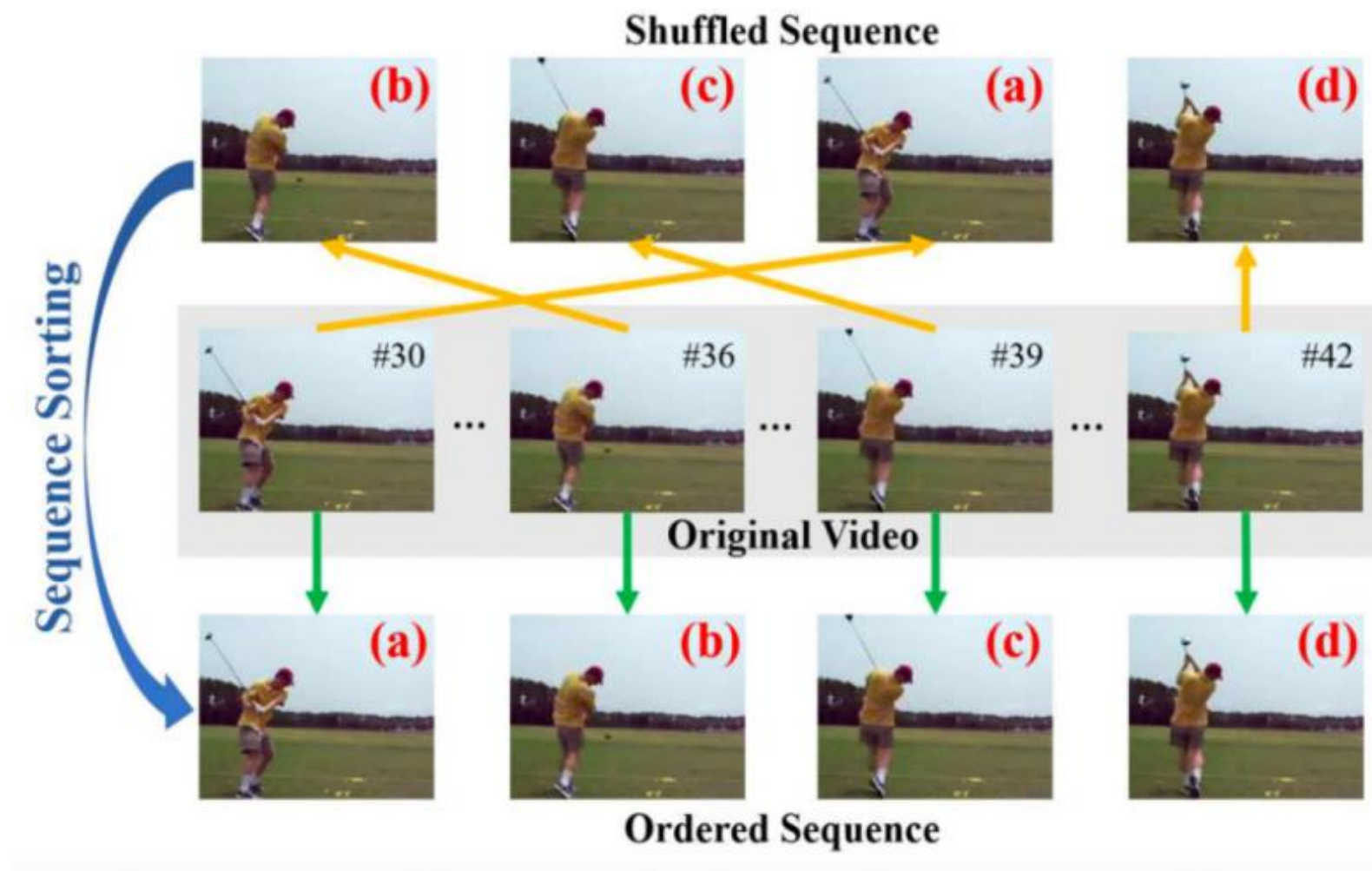
Main hypothesis of SSL

True or false?

The idea of SSL is that if the model can solve a complicated task that requires high-level understanding of the input, then it will learn more generalisable features

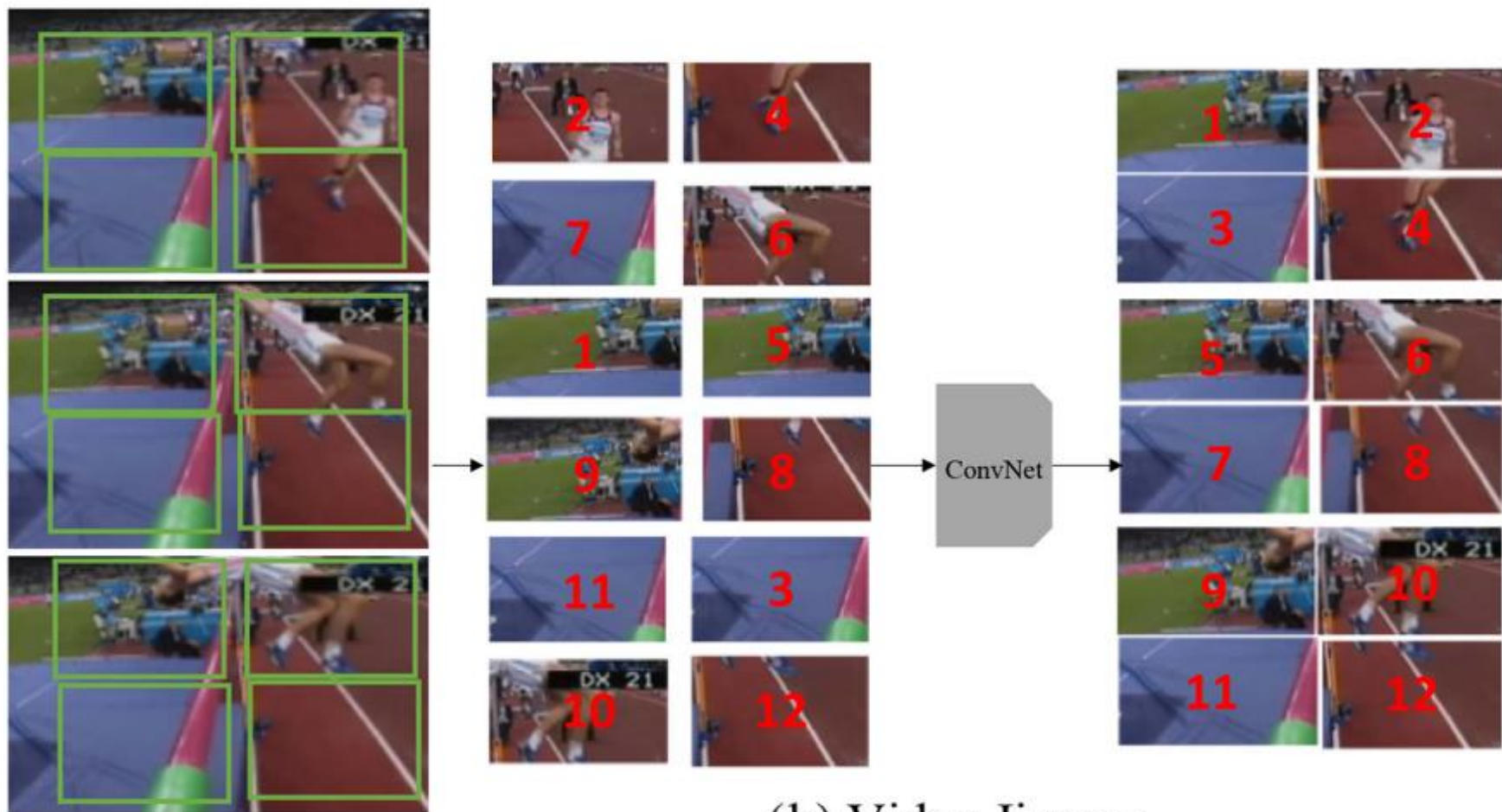


(a) Binary Classification Task



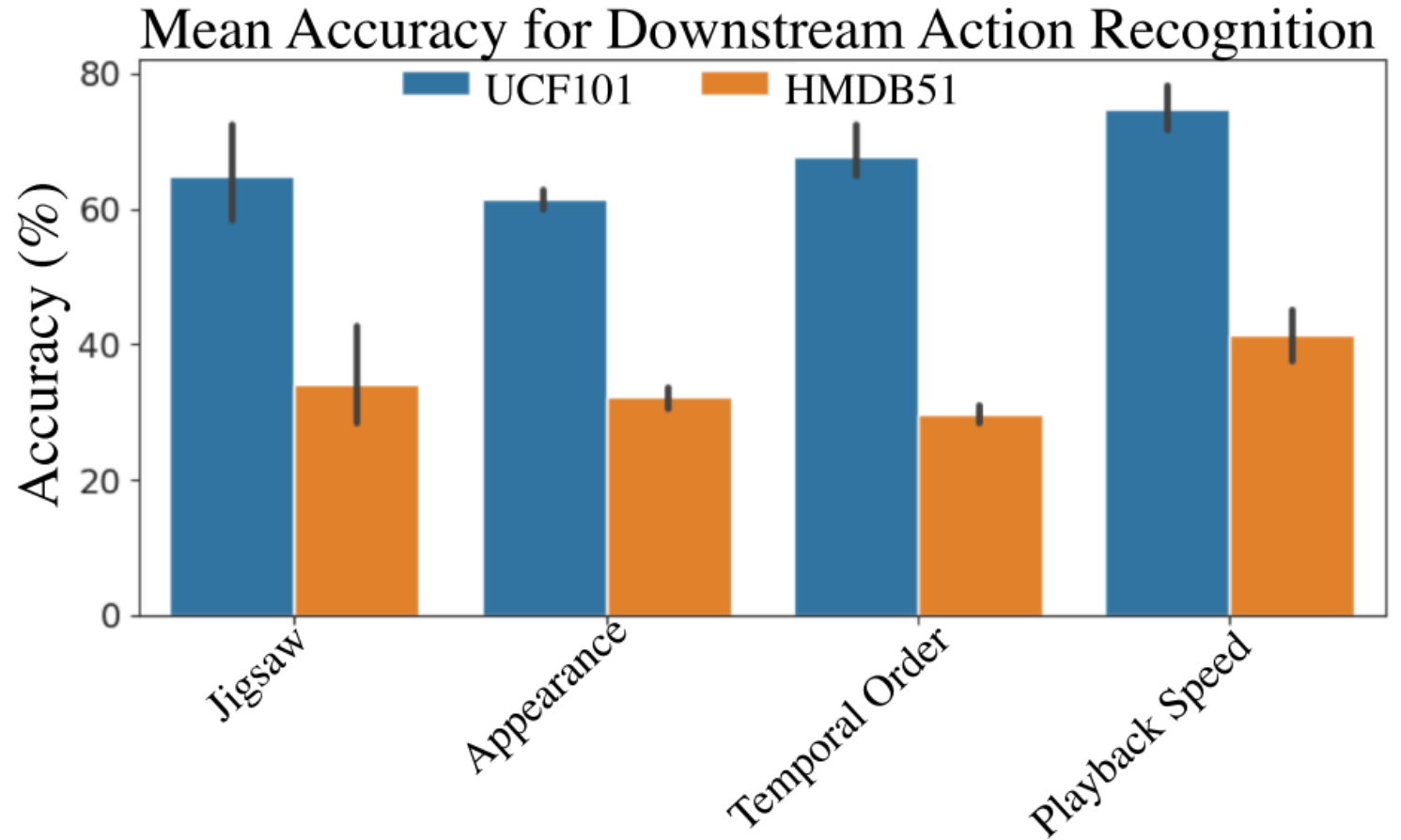
(b) Sequence Sorting





(b) Video Jigsaw

Do different pretext tasks make any difference?

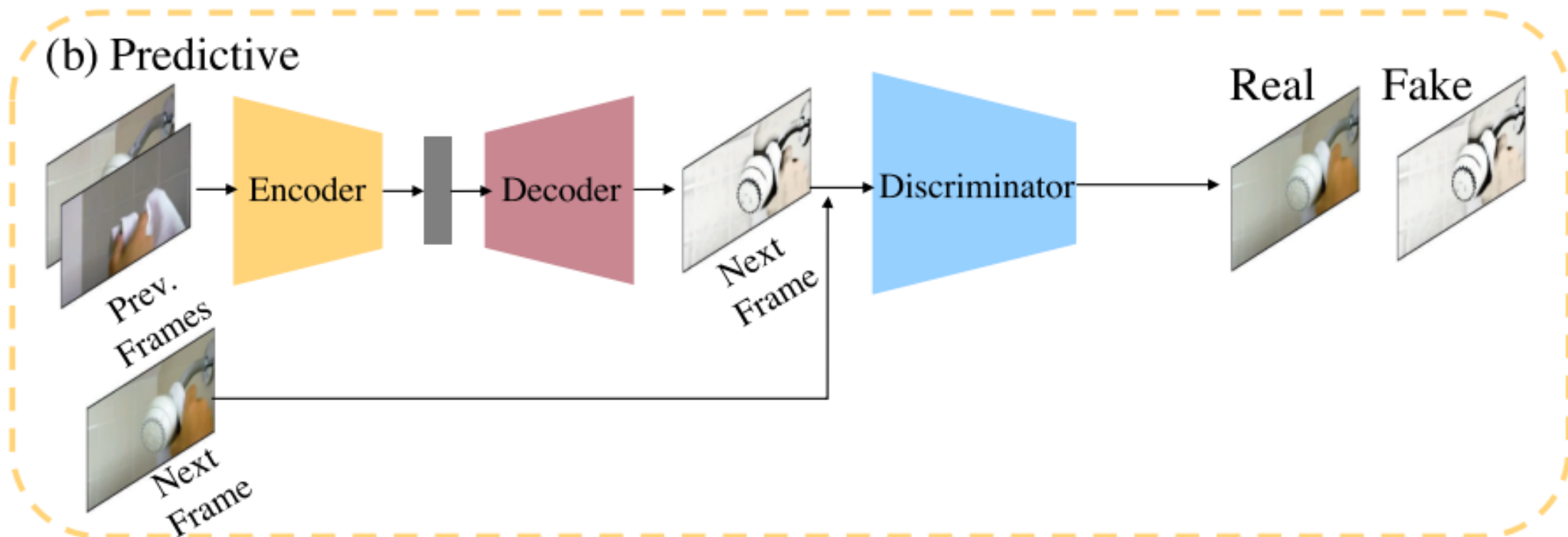
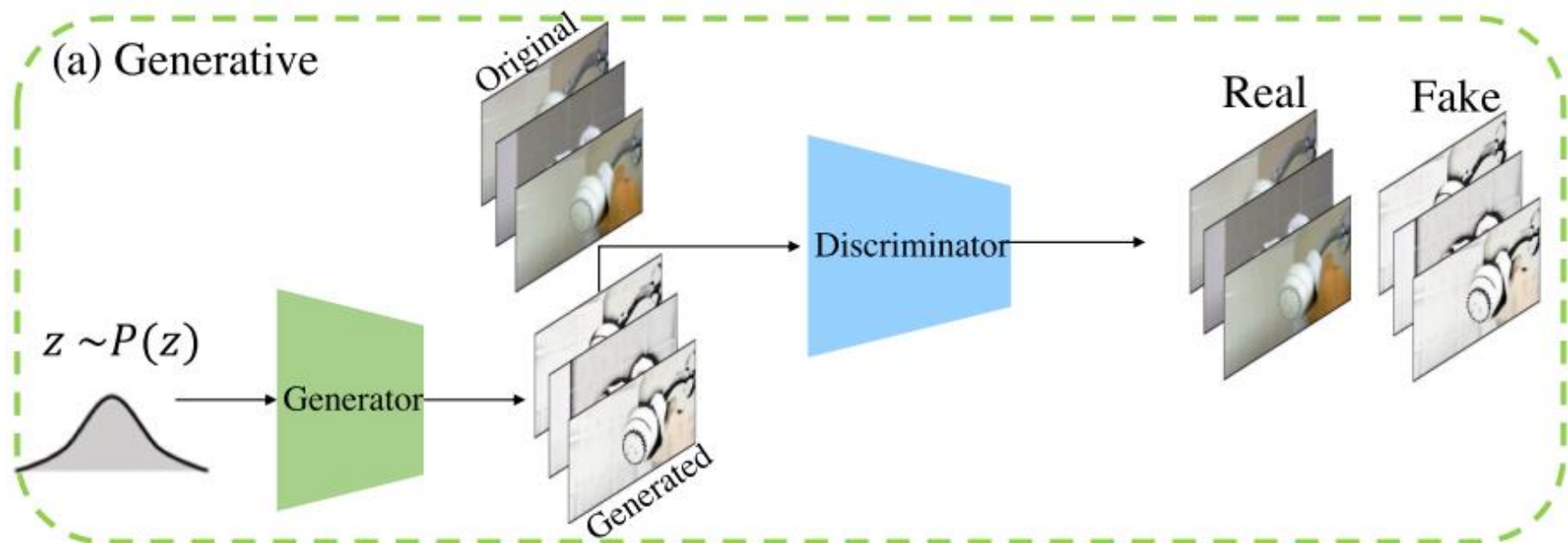


Are time-related pretext tasks more helpful?

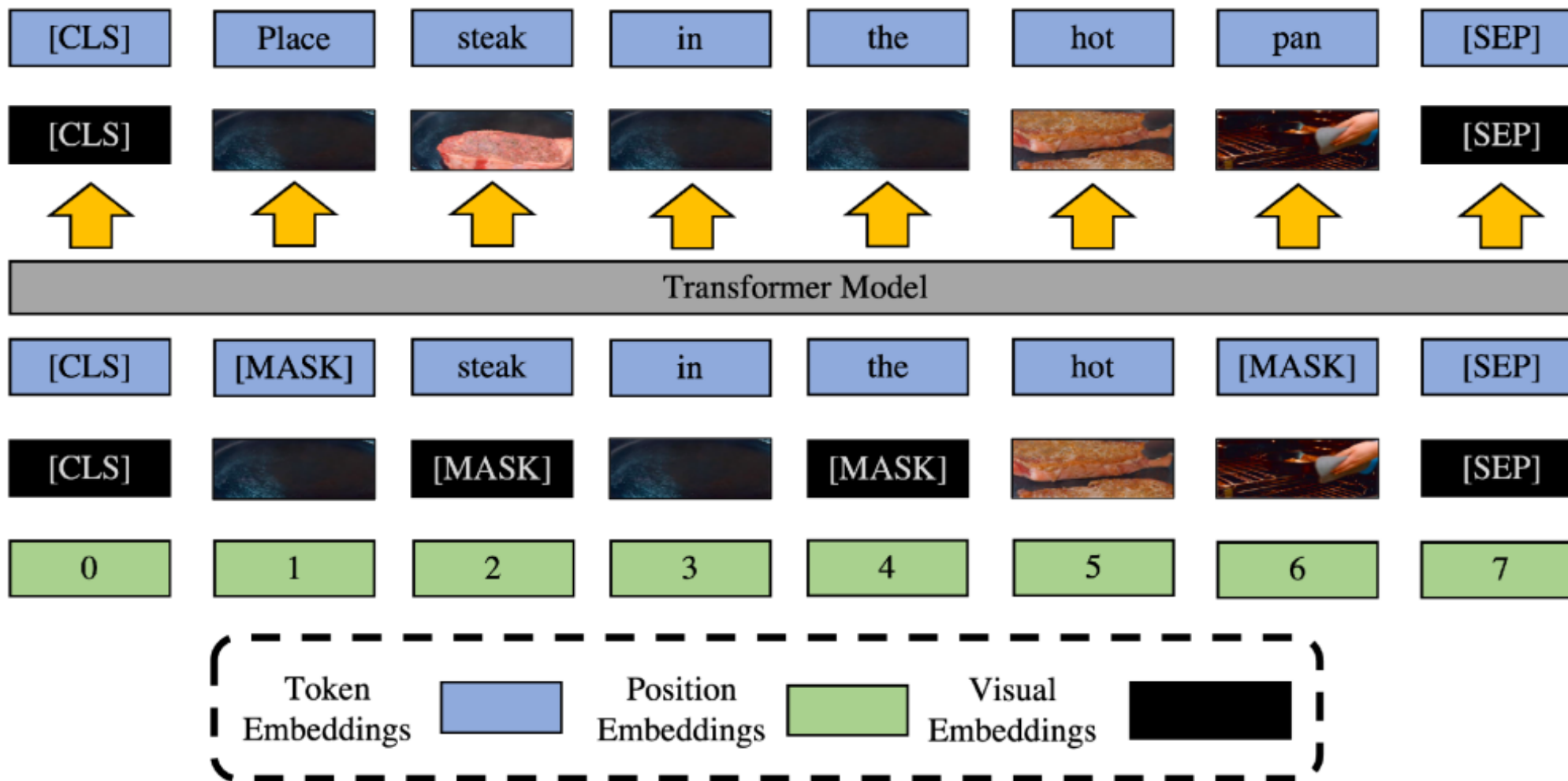
Generative approaches

GANs (Generative Adversarial Networks)
Masked Autoencoders (MAEs)

Predicting next frame



Masked modelling

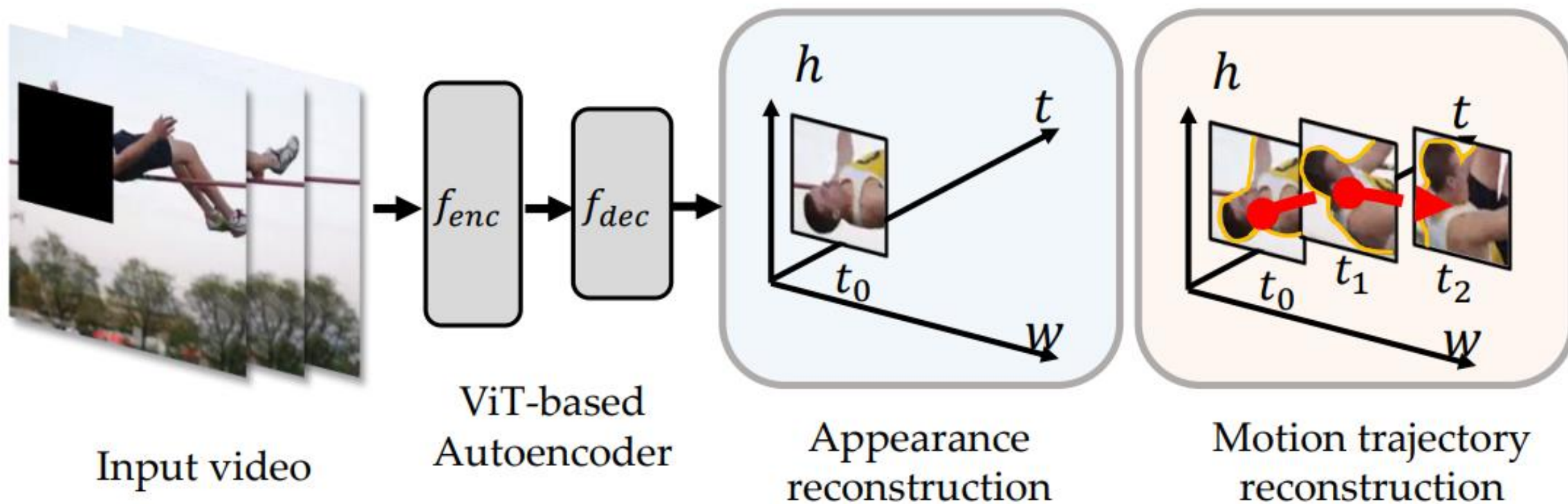




→ Position change

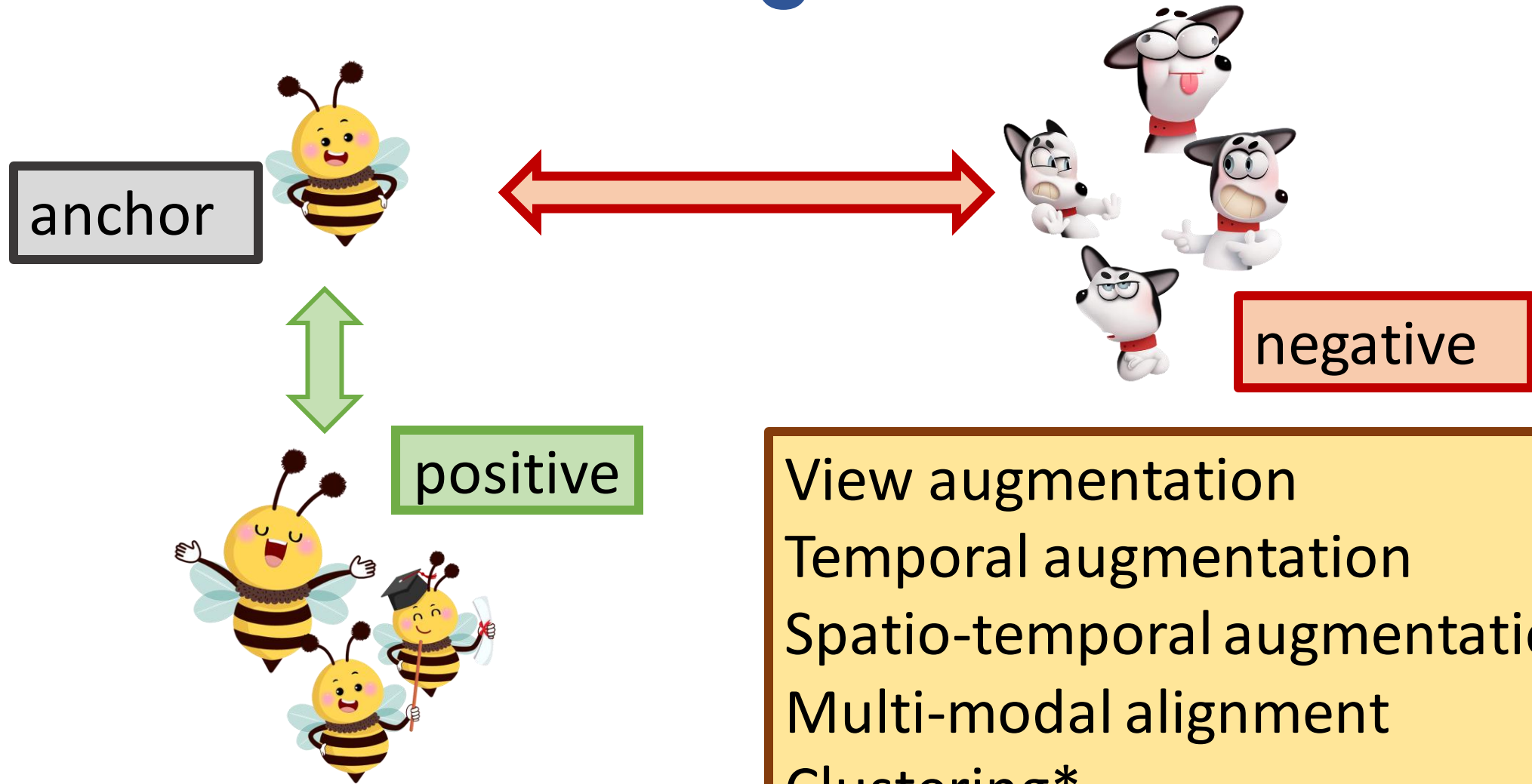
□ Shape change

(a) Two key factors to recognize a *high jump* action.



(b) Appearance reconstruction vs. motion trajectory reconstruction.

Contrastive learning



View augmentation
Temporal augmentation
Spatio-temporal augmentation
Multi-modal alignment
Clustering*

(*) Not seen here

Contrastive loss

$$\mathbf{L} = (1 - \mathbf{Y}) * ||\mathbf{x}_i - \mathbf{x}_j||^2 + \mathbf{Y} * \max(0, \mathbf{m} - ||\mathbf{x}_i - \mathbf{x}_j||^2)$$

$\mathbf{Y}=0$ if \mathbf{x}_i and \mathbf{x}_j have the same labels (and 1 otherwise)

Triplet loss

$$\mathbf{L} = \max(0, ||\mathbf{x} - \mathbf{x}^+||^2 - ||\mathbf{x} - \mathbf{x}^-||^2 + \mathbf{m})$$

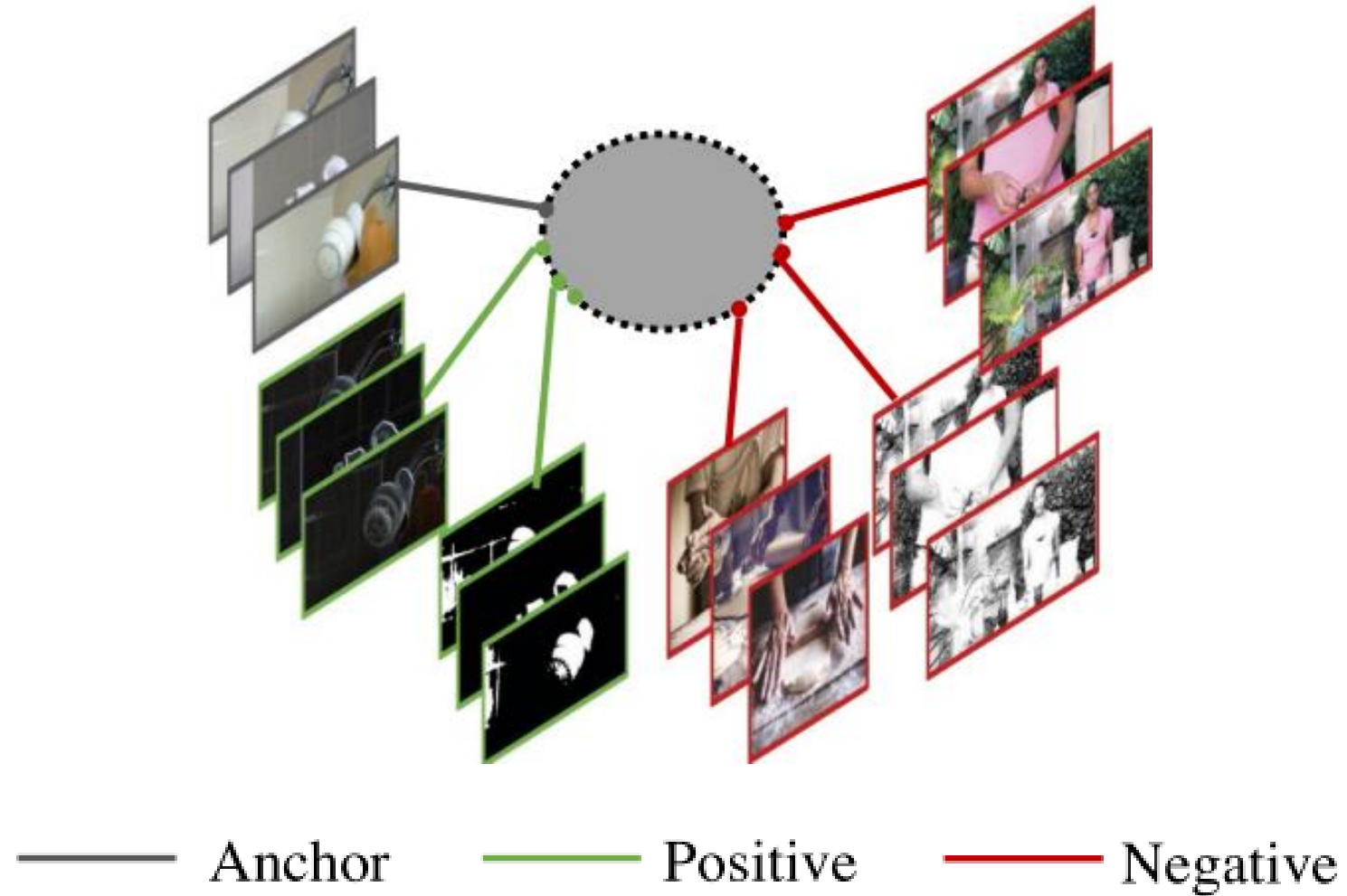
Do the contrastive loss and the triplet loss require the same data at the same time?

Noise Contrastive Estimation Loss (NCE)

Uses a pair of positive and a set of negative examples

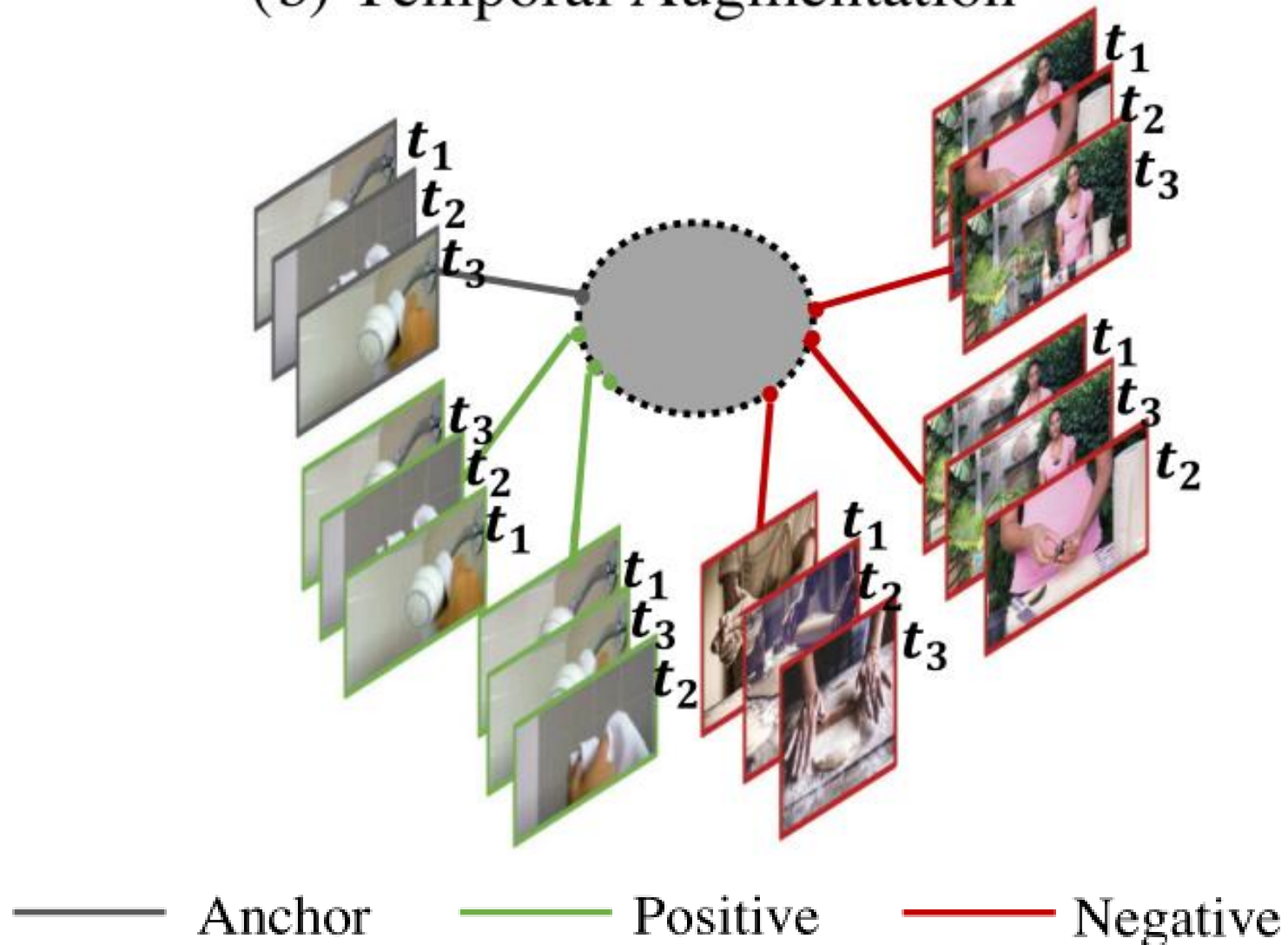
View augmentation

(a) View Augmentation



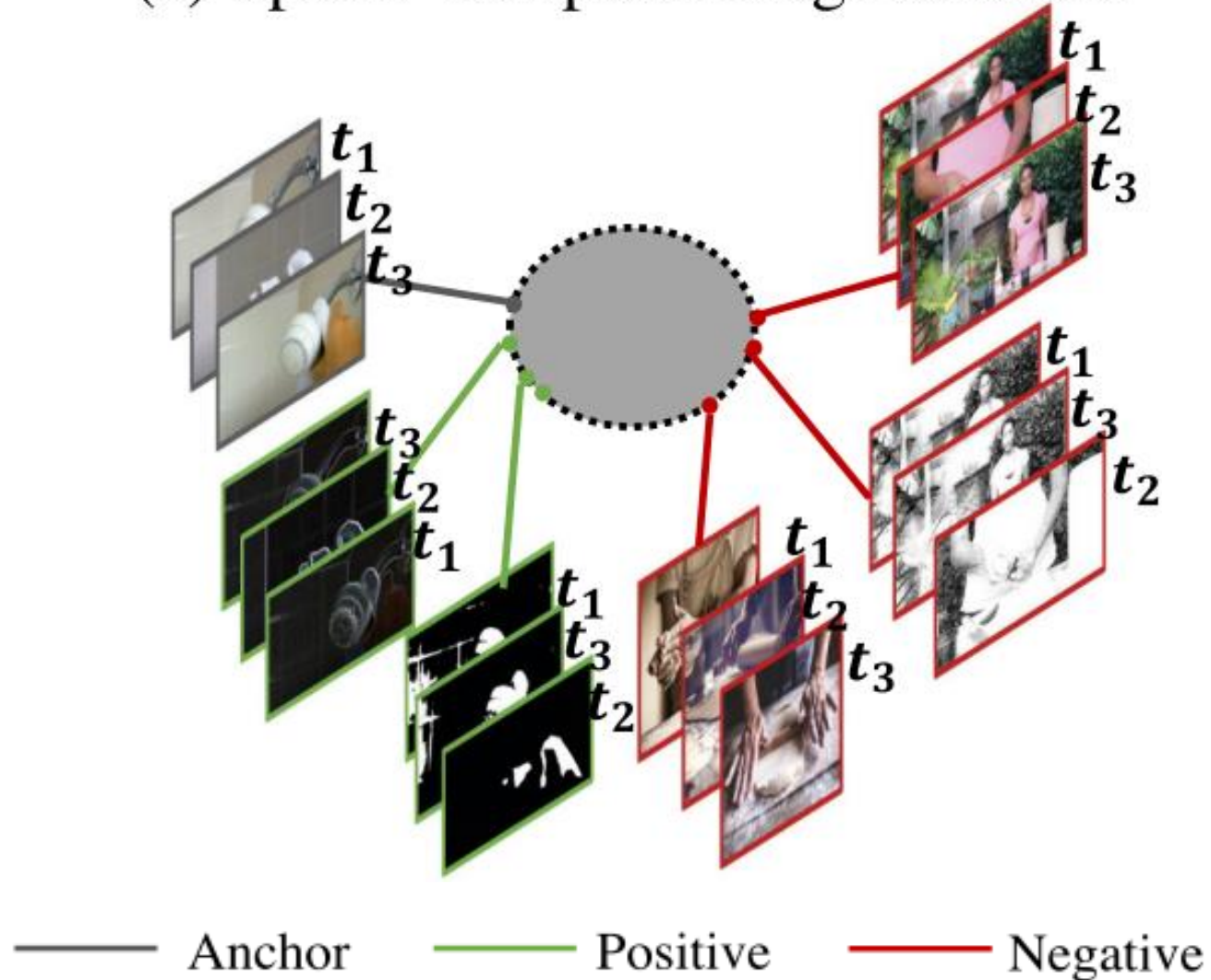
Temporal augmentation

(b) Temporal Augmentation



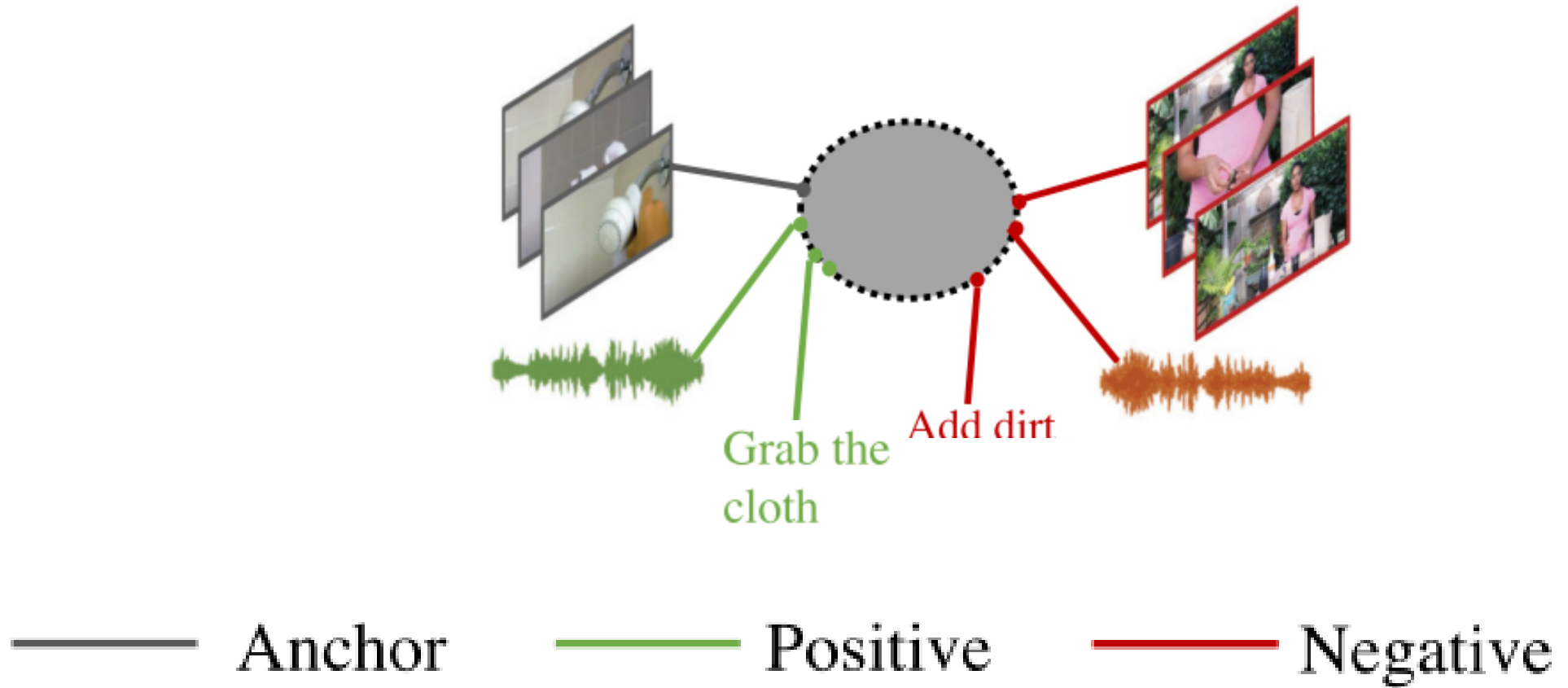
Spatio-temporal augmentation

(c) Spatio-Temporal Augmentation



Multimodal

(d) Cross-Modal Agreement



Comparison of contrastive approaches

