# Learning Object State Change in Videos: An Open-World Perspective

Zihui (Sherry) Xue[1,2]  Kumar Ashutosh[1,2]  Kristen Grauman[1,2]

[1]UT Austin  [2]FAIR, Meta

See our website for data, code & qualitative videos →

## Video Object State Change (OSC)

Objective: temporally localize an object's three states (initial, transitioning and end) from a video



## Motivation

- OSCs naturally exhibit a **long-tail**. Certain OSCs, such as melting butter, are frequently shown in videos while others like melting jaggery might be rarely seen.
- Prior works assume a **closed** vocabulary, limited to identifying state changes for objects observed during training.



Frequent OSCs

butter unmelted → melted

marshmallow unmelted → melted

Rare OSCs

novel object ?

→ **Open-world formulation**

## Quantitative Results

Quantitative results on ChangeIt [3] open-world (top) and HowToChange (bottom)



State Prec.@1 (known OSCs / novel OSCs):
CLIP [5] 0.29 / 0.29; VideoCLIP [6] 0.25 / 0.24; InternVideo [7] 0.29 / 0.25; LookFotheChange [3] 0.36 / 0.25; MultiTaskChange [4] 0.41 / 0.22; VIDOSC (ours) 0.56 / 0.48

F1 Score (%) (known OSCs / novel OSCs):
CLIP [5] 26.9 / 25.4; VideoCLIP [6] 36.6 / 34.3; InternVideo [7] 29.9 / 29.5; LookFotheChange [5] 30.3 / 28.7; MultiTaskChange [4] 33.9 / 29.9; VIDOSC (ours) 46.4 / 43.1

## Qualitative Results



Initial State | Transitioning State | End State

tying tie

peeling dragon fruit

grating orange

Top-1 frame predictions (known OSCs)

Initial State | Transitioning State | End State

tying rope

peeling avocado

grating cauliflower

Top-1 frame predictions (novel OSCs)

## Framework Overview



(a) Mining for OSC examples

[ASR transcription] you're going to use some rotisserie chicken so just get your rotisserie chicken and shred it up

LLAMA2

This video may contain the OSC of *chicken* + *shredding* (object + state transition)

(b) Pseudo Label Generation

State Description [whole chicken, shredding chicken, shredded chicken]

CLIP

State = Initial | State = Transitioning | ... | State = Background | State = End

Decoder / Video Encoder

$\tilde{\mathbf{y}}_1$ $\tilde{\mathbf{y}}_2$ $\tilde{\mathbf{y}}_3$ ... $\tilde{\mathbf{y}}_t$ $\tilde{\mathbf{y}}_{t+1}$ ... $\tilde{\mathbf{y}}_T$

$\mathbf{x}_1$ $\mathbf{x}_2$ $\mathbf{x}_3$ ... $\mathbf{x}_t$ $\mathbf{x}_{t+1}$ $\mathbf{x}_T$

□ object features

(d) Model Testing

Known OSCs: *shredding chicken* ... *shredding cabbage*

Novel OSCs: *shredding onion* ... *shredding coconut*

(c) Model Training

(a) Leverage ASR transcriptions and LLM for automatic OSC mining from instructional videos;
(b) Employ OSC textual descriptions with VLM for pseudo label generation;
(c) Video OSC model for object-agnostic state predictions (shared state vocabulary, temporal modeling, object-centric features);
(d) Test with open-world formulation, evaluating performance on both known and novel OSCs.

## HowToChange Dataset

The first open-world benchmark for video OSC localization



| | #Obj | #ST | #OSC | #Obj per ST | #Video | GT Label? |
|---|---|---|---|---|---|---|
| Alayrac et al. [1] | 5 | 6 | 7 | 1.2 | 630 | ✅ |
| TaskFluent [2] | 25 | 14 | 32 | 2.3 | 809 | ✅ |
| ChangeIt (Train) [3] | 42 | 27 | 44 | 1.6 | 34,428 | |
| ChangeIt (Eval) [3] | 42 | 27 | 44 | 1.6 | 667 | ✅ |
| HowToChange (Train) | 122 | 20 | 318 | 15.9 | 36,075 | |
| HowToChange (Eval) | 134 | 20 | 409 | 20.5 | 5,424 | ✅ |

Obj = object, ST = state transition

Features:
- Scale jump in the number of OSCs (#OSC) and videos (#Video)
- Wide variety of objects per state transition (#Obj per ST)

Ground truth annotation distribution



Annotation Distribution by State Transition Type (known / novel)

Annotation Distribution by Object Type (known / novel)

## Qualitative Results

Comparison of model predictions on one test video of slicing shallot



CLIP [5]
VideoCLIP [6]
InternVideo [7]
LookFotheChange [3]
MultiTaskChange [4]
VIDOSC (ours)
Ground Truth

Background | Initial State | Transitioning State | End State

[1] Alayrac et al., Joint discovery of object states and manipulation actions, ICCV 17.
[2] Liu et al., Jointly recognizing object fluents and tasks in egocentric videos, ICCV 17.
[3] Soucek et al., Look for the Change: Learning Object States and State-Modifying Actions from Untrimmed Web Videos, CVPR 22.
[4] Soucek et al., Multi-Task Learning of Object State Changes from Uncurated Videos, arXiv 22.
[5] Radford et al., Learning transferable visual models from natural language supervision, ICML 21.
[6] Xu et al., VideoCLIP: Contrastive Pre-training for Zero-shot Video-Text Understanding, EMNLP 21.
[7] Wang et al., InternVideo: General Video Foundation Models via Generative and Discriminative Learning, arXiv 22.