

Challenges

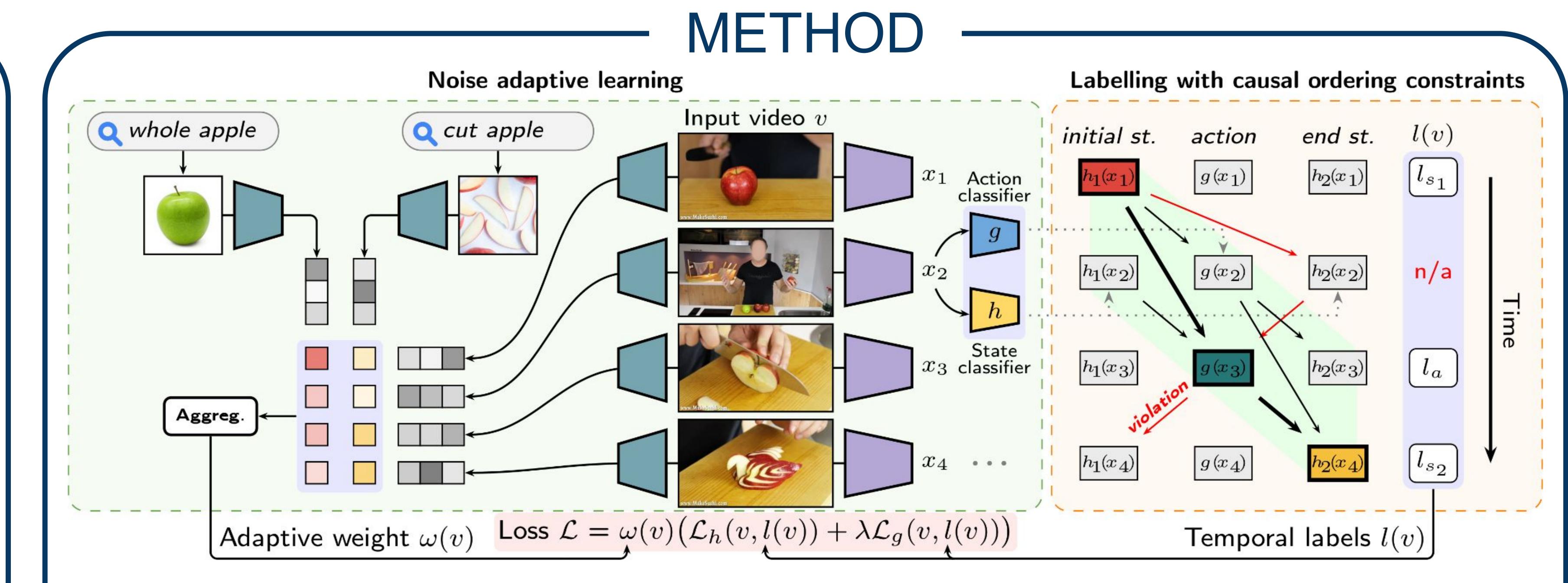
- > Visual variability of the objects and its states.
- \succ Thousands of objects with many more states, annotating is both difficult and expensive.
- In-the-wild, noisy, uncurated, and long videos.

Contributions

- Self-supervised model for learning object states and state-modifying actions from long uncurated web videos.
- > Causal ordering signal (*initial state* \rightarrow *action* \rightarrow *end state*) is used as the supervision.
- New uncurated dataset with 2600+ hours of video and 34 thousand changes of object states.

Look for the Change: Learning Object States and **State-Modifying Actions from Untrimmed Web Videos**

Tomáš Souček¹, Jean-Baptiste Alayrac², Antoine Miech², Ivan Laptev³, Josef Sivic¹ ¹Czech Technical University ²DeepMind ³ENS/INRIA



Learning step

cut an apple

 \succ compute loss \mathcal{L} and perform a gradient update on the classifiers using the computed labels l(v)

Noise adaptive weighting $\omega(v)$

training videos may be unrelated \rightarrow we compute similarity between video frames and example images $\mathcal{E}_1, \mathcal{E}_2$ and use normalized similarity score $\omega(v) = \sigma\left(\frac{r_v - \theta}{\tau}\right)$ to weight contribution of the video in the loss

cut a pineapple

"How to..

fold a paper plane?"

 $r_v = \max_{t < t'} \sum_{t < t'} \sin(e_1, v_t) \sum_{t < t'} \sin(e_2, v_{t'})$

Labelling step

Computing Labels l(v)

labels l(v) is a triplet of the *initial state* l_{s_1} , action l_a and end state l_{s_2} positions.

The labels satisfy the *causal ordering constraint*: $\mathcal{D}_{v} = \{ l : 1 \le l_{s_{1}} < l_{a} < l_{s_{2}} \le T_{v} \}$

CHANGEIT DATASE1 "How to..

- "How to cut an apple?"
- cutting, ball inflating, etc.

> compute labels l(v) for a video $v = \{x_t\}_{t=1}^{T_v}$ using the state h and action g classifiers:

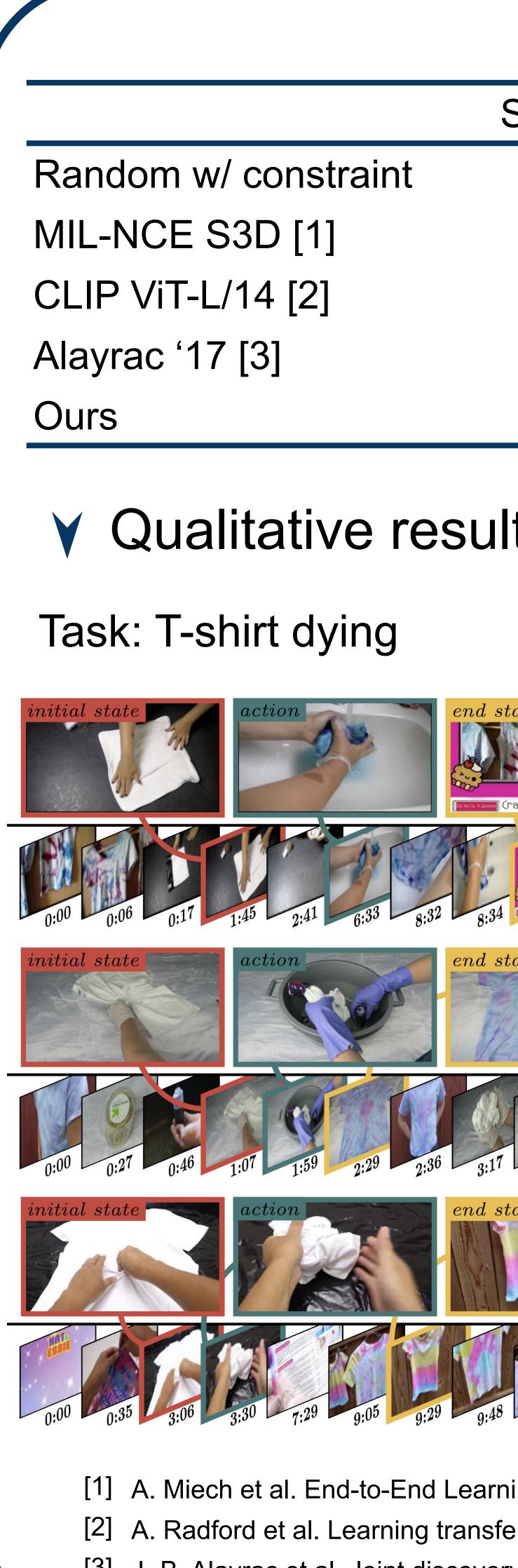
 $l(v) = \arg\max h_1(x_{l_{s_1}}) \cdot g(x_{l_a}) \cdot h_2(x_{l_{s_2}})$

Gathered by searching Youtube for terms such as

> 44 state-changing action categories such as apple

> 34,428 in-the-wild videos, in total 2,642 hours, average video length 4.6 minutes

> 667 videos per-frame annotated with labels: background, initial state, action, end state



	RESULTS -	
St prec Ac prec Random w/ constraint 0.15 0.41 MIL-NCE S3D [1] 0.27 0.50 CLIP VIT-L/14 [2] 0.30 0.63 Alayrac '17 [3] 0.30 0.59 Ours 0.35 0.68	 dataset A separate mode A single frame pr 	ategories of the Changelt el trained for each category redicted as the initial state le end state per video
V Qualitative results, see the co	onsistency	
Task: T-shirt dying Task:	: Cream whisking	Task: Paper plane folding
initial state initial state initial state initial state initial state initial state	SUBSCRIBE NOW action SUBSCRIBE NOW Image: state Image: state Image: state Image: state Image: state Image:	initial state <code>n the market</code>
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	0:12 0:36 1:43 5:30 6:22 6:23 6:24 6:43	$ \underbrace{ \begin{array}{c} \hline \\ \hline $
initial state	action end state	initial state Set aside Design and state Design and state Set aside Design and state Set aside Set
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	$ \begin{array}{c} \begin{array}{c} \begin{array}{c} \begin{array}{c} \\ \\ \\ \\ \end{array} \end{array} \end{array} \\ \begin{array}{c} \\ \\ \\ \\ \\ \\ \\ \\ \\ \\ \\ \\ \\ \\ \\ \\ \\ \\ \\$
initial state	action end state for the state fo	initial state
$ \begin{array}{c} \hline & & & & \\ \hline & & & \\ 0:00 & & & \\ 0:35 & & & \\ 0:35 & & & \\ 3:06 & & \\ \hline & & \\ 3:06 & & \\ \hline & & \\ 7:29 & & \\ 9:05 & & \\ 9:05 & & \\ 9:29 & & \\ 9:29 & & \\ 9:48 & & \\ 10:11 & \\ 10:48 & \\ 10:11 & \\ 10:48 & \\ \hline & \\ 10:48 & \\ \hline & \\ 0:00 & \\ 0:06 & \\ \hline \end{array} \right) $	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$
 A. Miech et al. End-to-End Learning of Visual Representatio A. Radford et al. Learning transferable visual models from n JB. Alayrac et al. Joint discovery of object states and mani 	atural language supervision, 2021.	VPR, 2019.
ABI A7	TIONS	
St prec Ac prec university w/o noise adapt. 0.34 0.64 only a single layer 0.34 0.61 w/o data augment. 0.33 0.66 Ours 0.35 0.68	3.4k 7k 17k 34k 3.4k 7 Number of training v	videos

