

Learning State-Aware Visual Representations from Audible Interactions



Himangi Mittal



Pedro Morgado



Unnat Jain



Abhinav Gupta

Neural Information Processing Systems (NeurIPS), 2022.

Real-world videos consists of long, untrimmed, egocentric daily activities



Real-world videos consists of long, untrimmed, egocentric daily activities



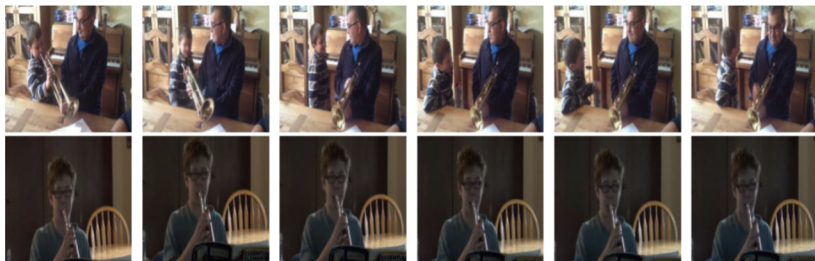
Learning representations from real-world videos can be quite challenging

Traditional Pipelines (1)

Work on curated, trimmed videos

Traditional Pipelines (1)

Work on curated, trimmed videos



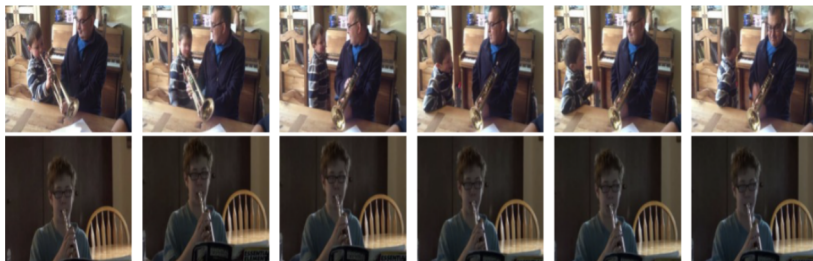
Playing trumpet



Dribbling basketball

Traditional Pipelines (1)

Work on curated, trimmed videos



Playing trumpet



Dribbling basketball

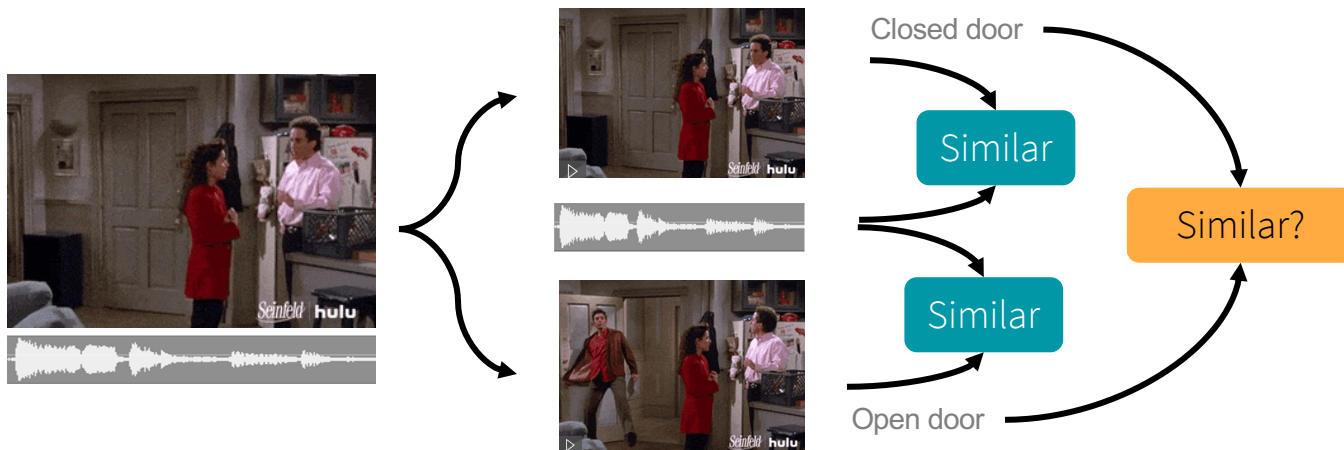
Long, untrimmed, real-world videos consists of multiple actions as well as **no-activity segments**

Traditional Pipelines (2)



Traditional Pipelines (2)

Invariant to state-changes in the environment



Can we learn meaningful representations from interaction-rich, untrimmed, and multi-modal egocentric data?

Two key components/contributions:

Focus learning on
moments of interaction

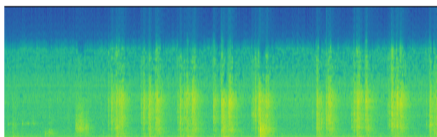
Learning from audible
state changes

Video



Cut celery

Spectrogram

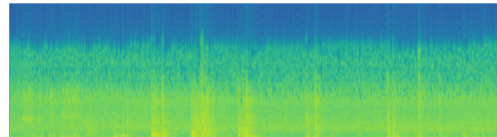


Video



Peel onion

Spectrogram

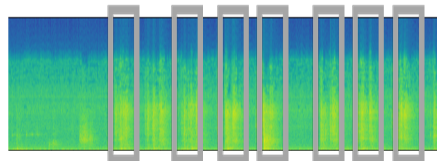


Video



Cut celery

Spectrogram

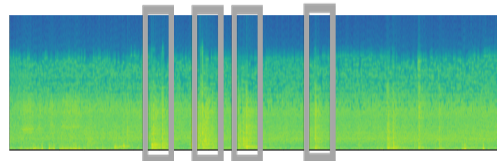


Video



Peel onion

Spectrogram

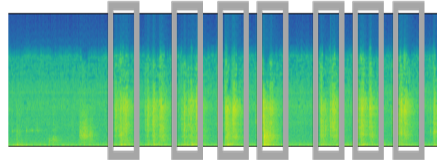


Video



Cut celery

Spectrogram

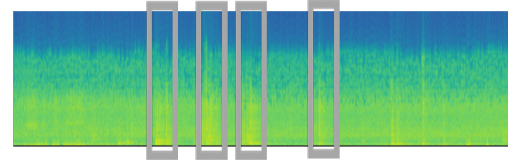


Video



Peel onion

Spectrogram

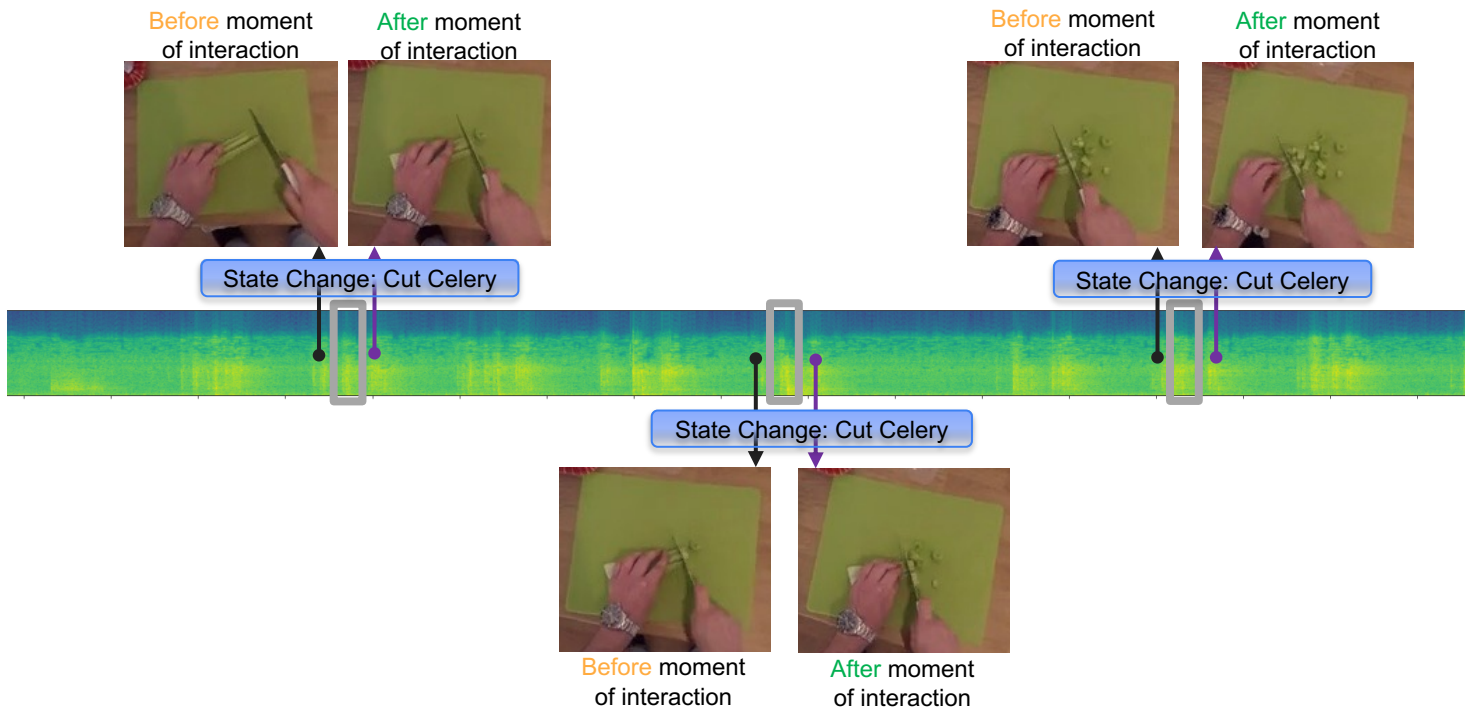


We refer these as **Moments of Interaction (Mol)**

Two key components/contributions:

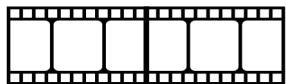
Focus learning on
moments of interaction

Learning from audible
state changes



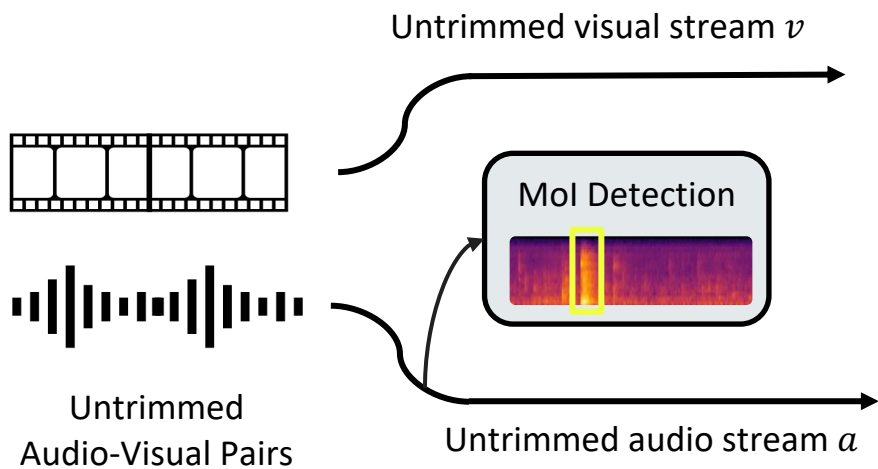
Transition between object states is often marked by characteristic sounds

Method

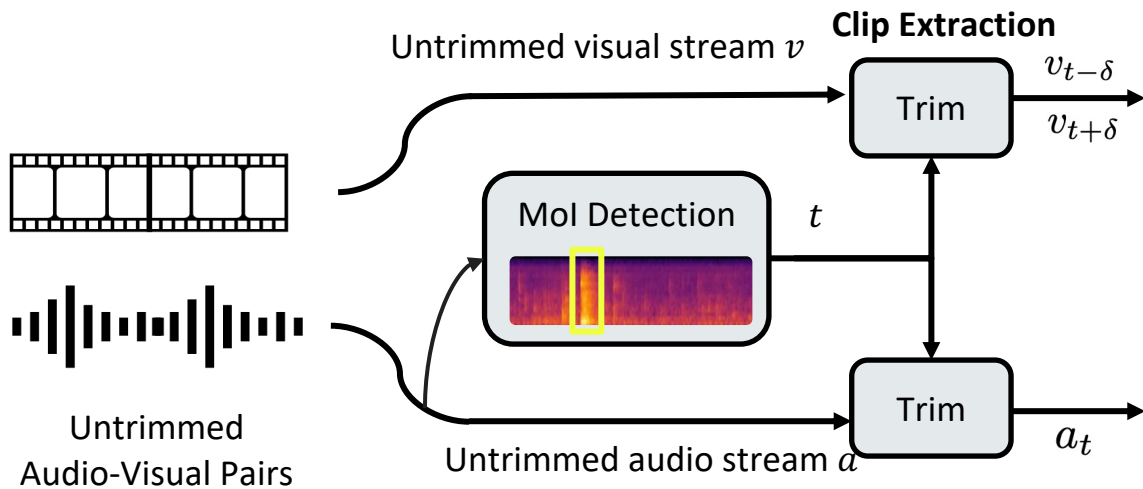


Untrimmed
Audio-Visual Pairs

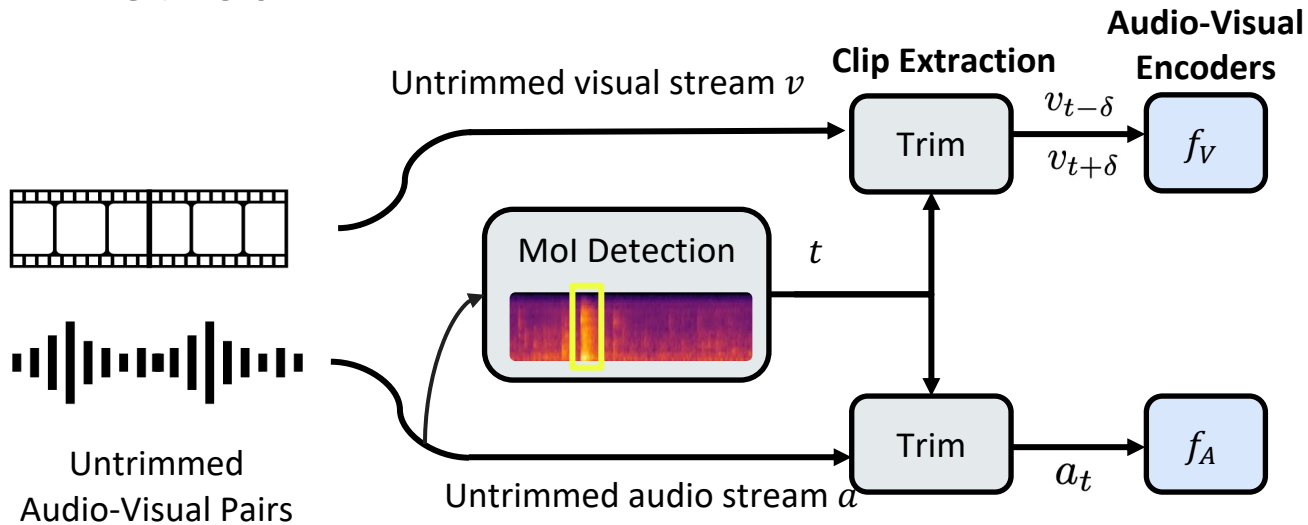
Method



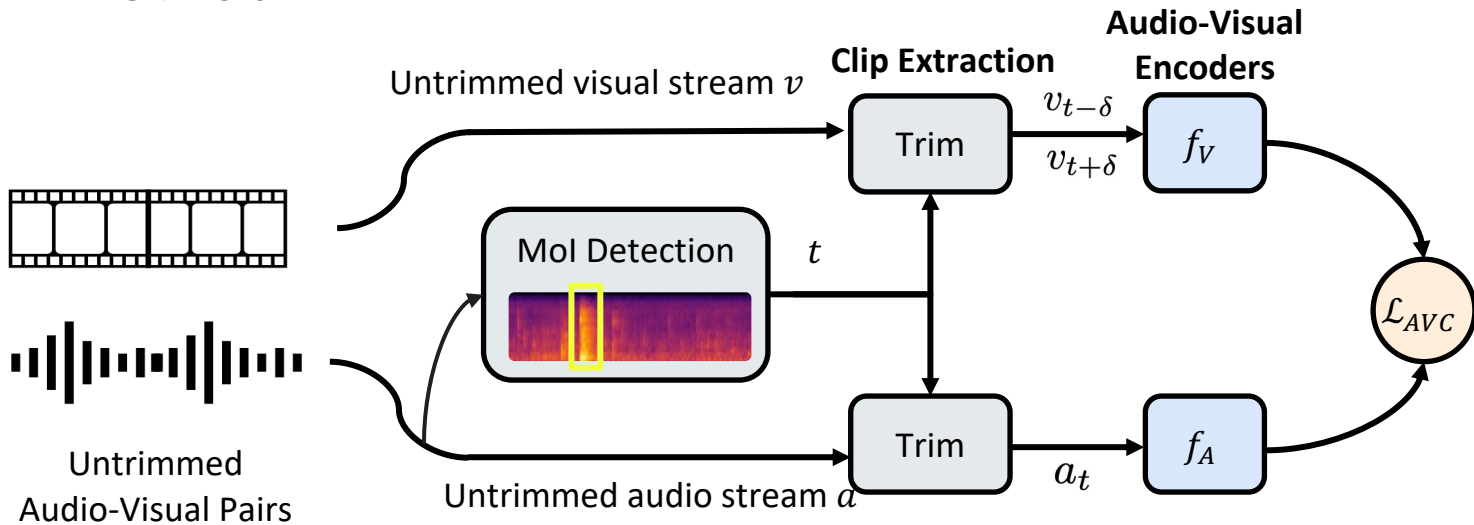
Method



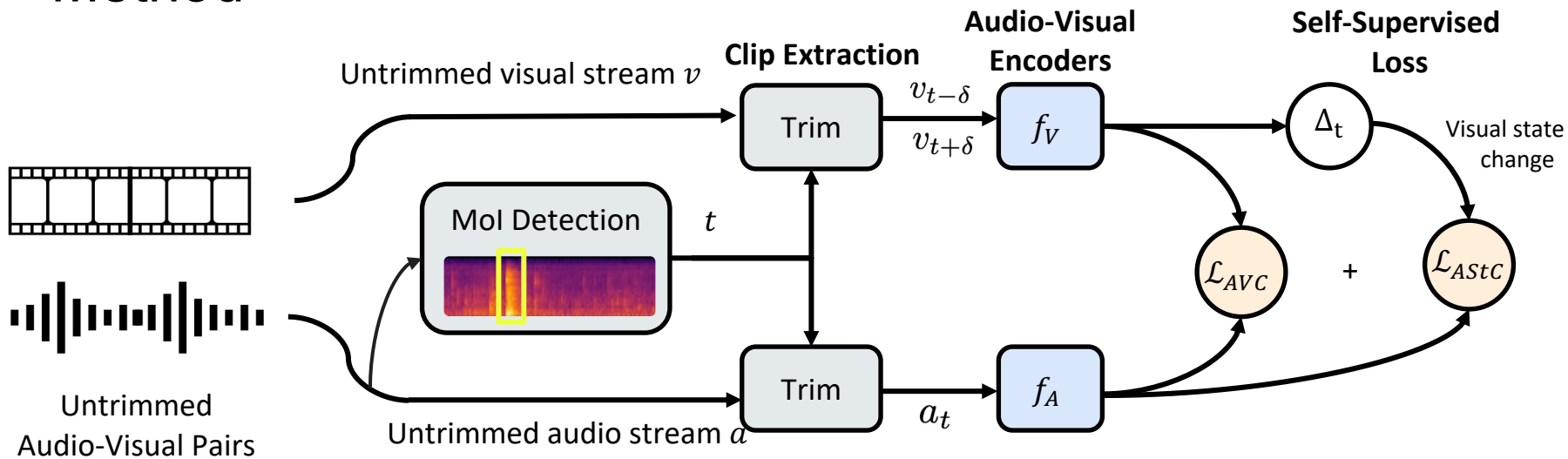
Method



Method

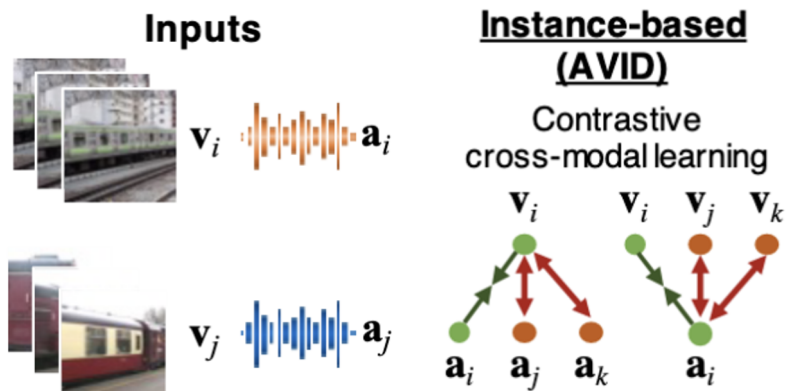


Method



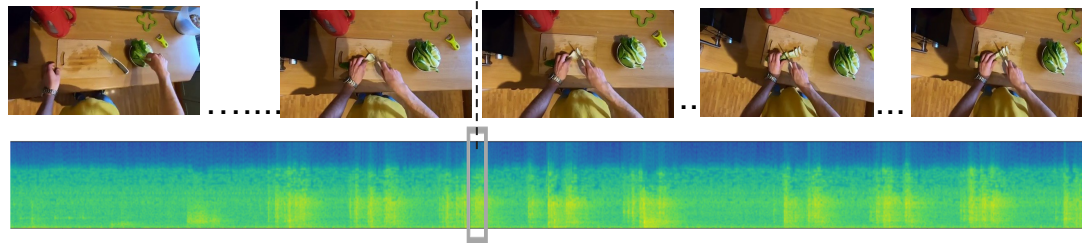
Audio-Visual Correspondence Loss (\mathcal{L}_{AVC})

Learns audio-visual representations by contrasting visual representations from multiple audios.



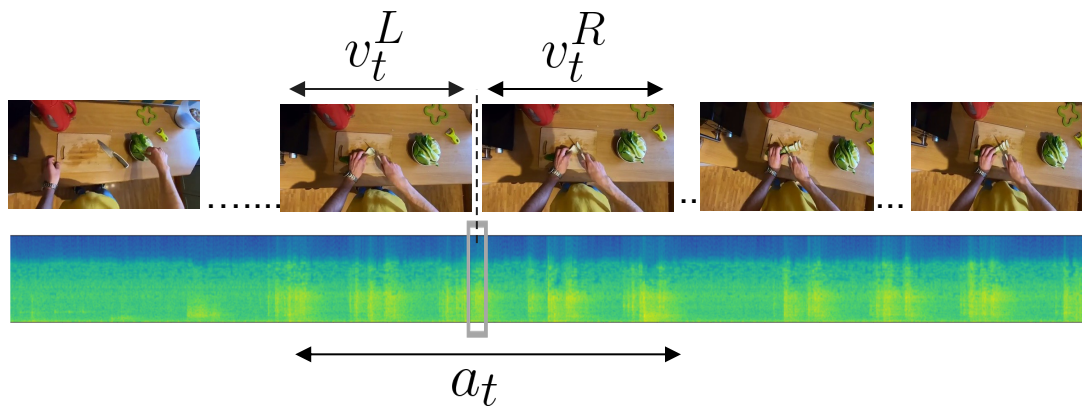
Audio-Visual Correspondence Loss (\mathcal{L}_{AVC})

Learns audio-visual representations by contrasting visual representations from multiple audios.



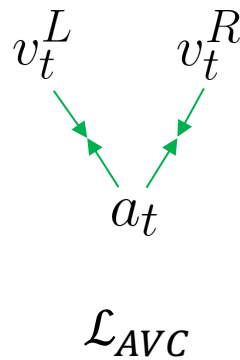
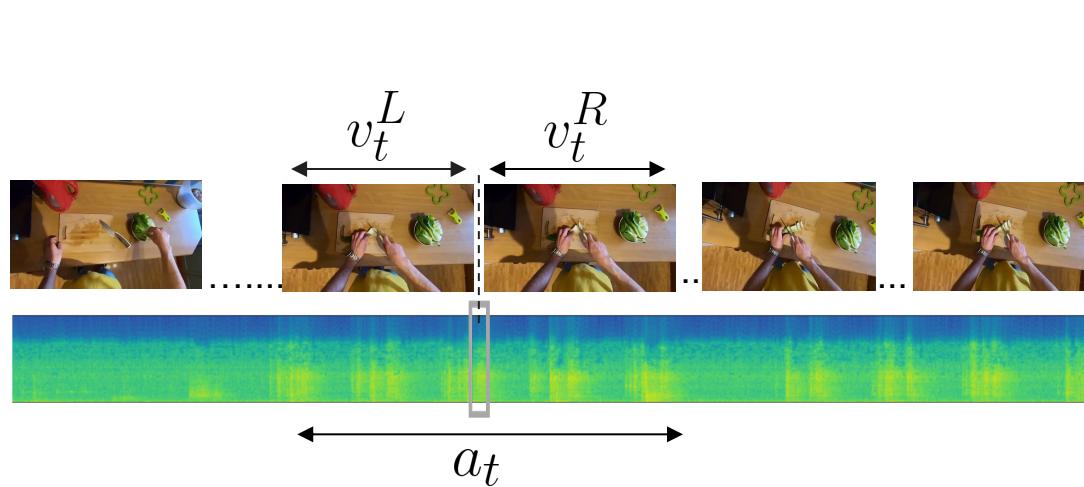
Audio-Visual Correspondence Loss (\mathcal{L}_{AVC})

Learns audio-visual representations by contrasting visual representations from multiple audios.



Audio-Visual Correspondence Loss (\mathcal{L}_{AVC})

Learns audio-visual representations by contrasting visual representations from multiple audios.

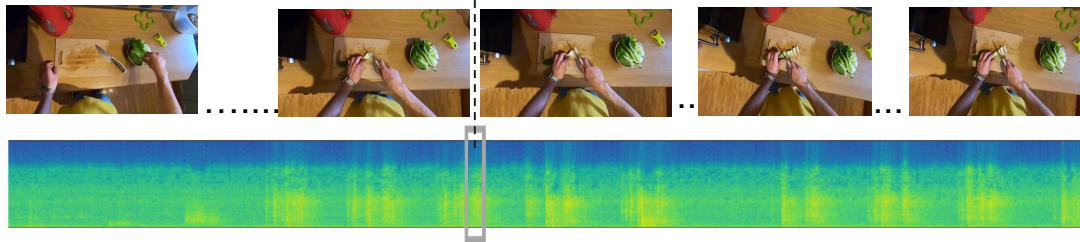


Audible State Change Loss (\mathcal{L}_{AstC})

Our proposed objective function tries to *associate the audio with changes in the visual state* during a moment of interaction.

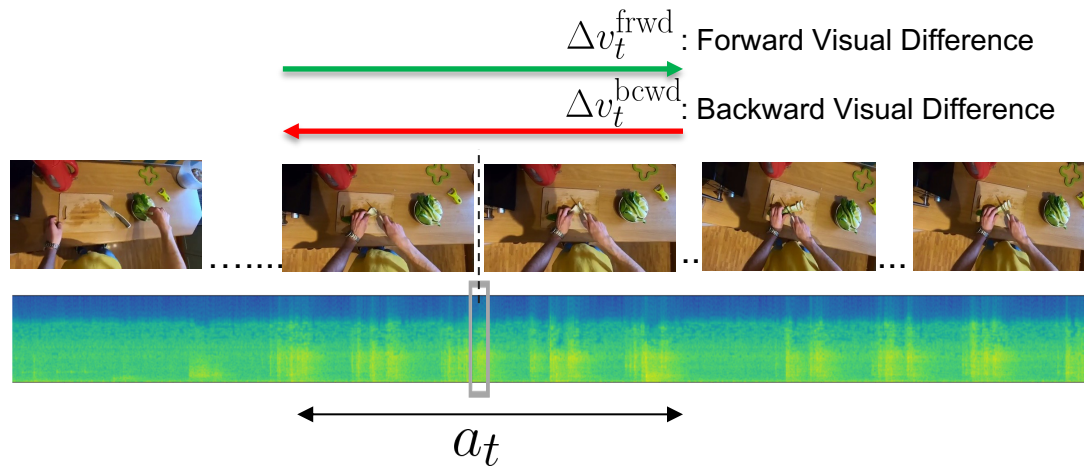
Audible State Change Loss (\mathcal{L}_{AstC})

Our proposed objective function tries to *associate the audio with changes in the visual state* during a moment of interaction.



Audible State Change Loss (\mathcal{L}_{AstC})

Our proposed objective function tries to *associate the audio with changes in the visual state during a moment of interaction.*

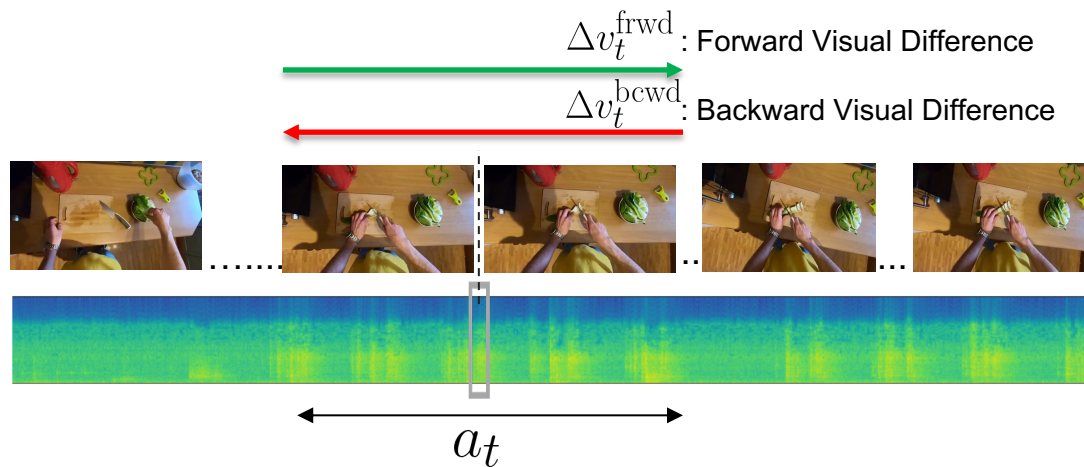


$$\Delta v_t^{frwd} = \phi(v_T^R - v_T^L)$$

$$\Delta v_t^{bcwd} = \phi(v_T^L - v_T^R)$$

Audible State Change Loss (\mathcal{L}_{AstC})

Our proposed objective function tries to *associate the audio with changes in the visual state during a moment of interaction.*



$$\Delta v_t^{\text{frwd}} = \phi(v_T^R - v_T^L)$$

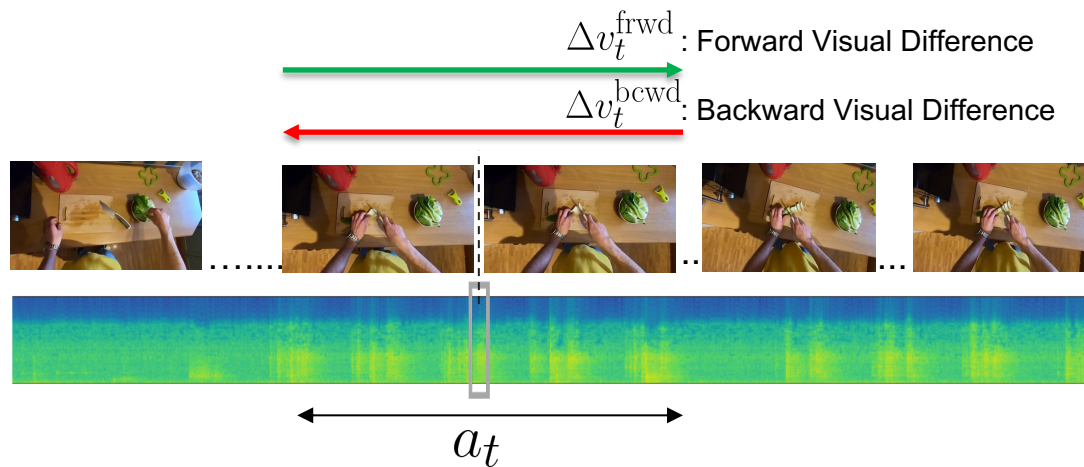
$$\Delta v_t^{\text{bcwd}} = \phi(v_T^L - v_T^R)$$

$$\langle \Delta v_t^{\text{frwd}}, a_t \rangle \uparrow$$

$$\mathcal{L}_{AstC}$$

Audible State Change Loss (\mathcal{L}_{AstC})

Our proposed objective function tries to *associate the audio with changes in the visual state during a moment of interaction.*



$$\Delta v_t^{frwd} = \phi(v_T^R - v_T^L)$$

$$\Delta v_t^{bcwd} = \phi(v_T^L - v_T^R)$$

$$\langle \Delta v_t^{frwd}, a_t \rangle \uparrow \quad \langle \Delta v_t^{bcwd}, a_t \rangle \downarrow$$

$$\mathcal{L}_{AstC}$$

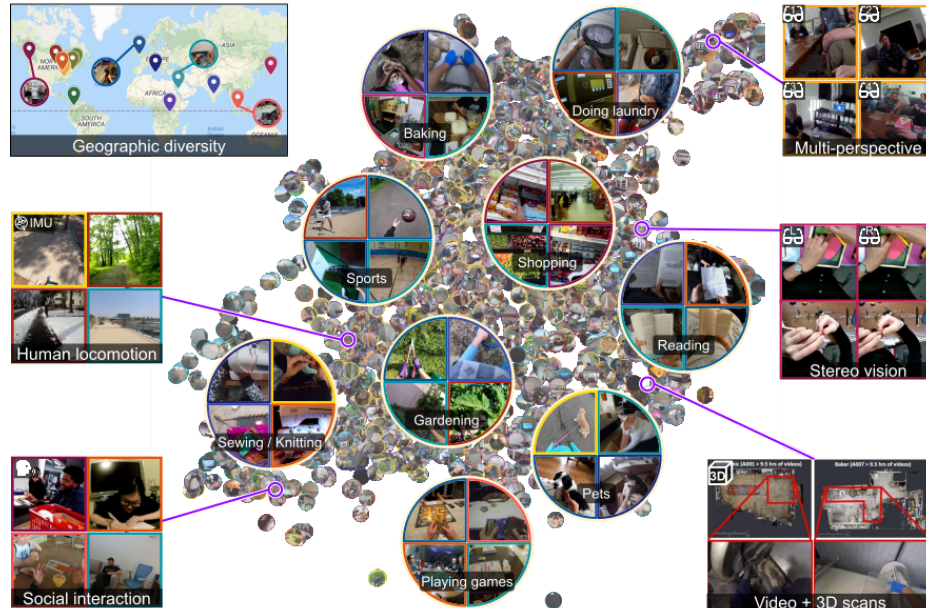
Datasets (1)

EPIC-Kitchens-100 : Consists of 100 hours of activities in the kitchen



Datasets (2)

Ego4D : Contains 3,670 hours of ego-centric video covering daily activities in the home, workplace, social settings, etc.



Quantitative Results on Action Recognition

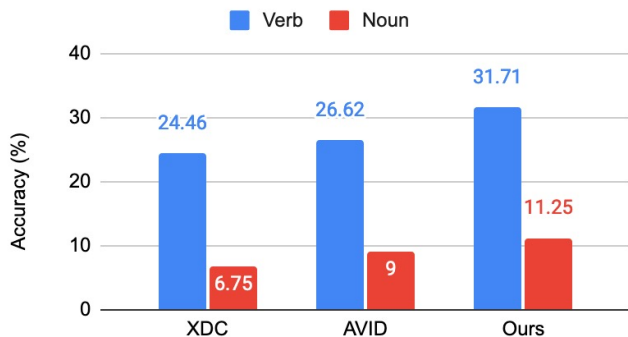


Predict the verb and noun of an action

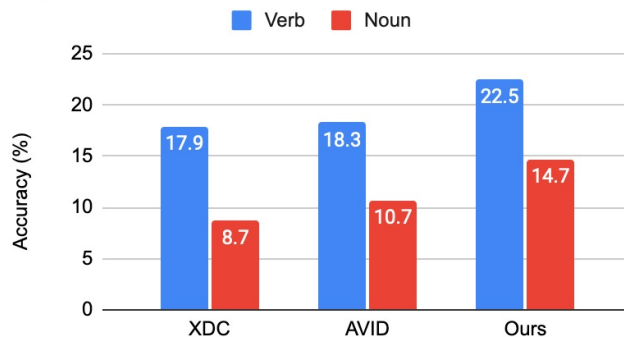
Discussion of Results (1)

Compared to the **baselines – XDC and AVID**, **our method** performs better by significant margins.

EPIC-Kitchens-100



Ego4D

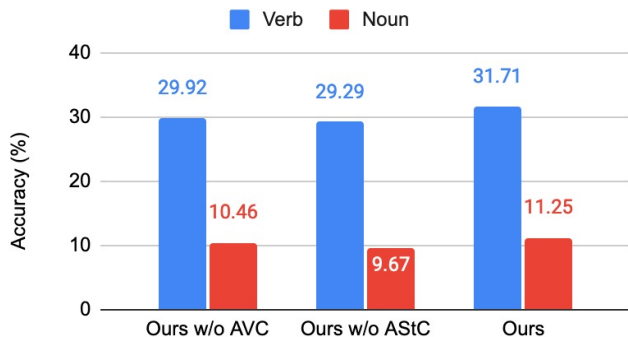


Action Recognition (Top-1 Accuracy) on EPIC-Kitchens-100 (left) and Ego4D (right)
Higher is better.

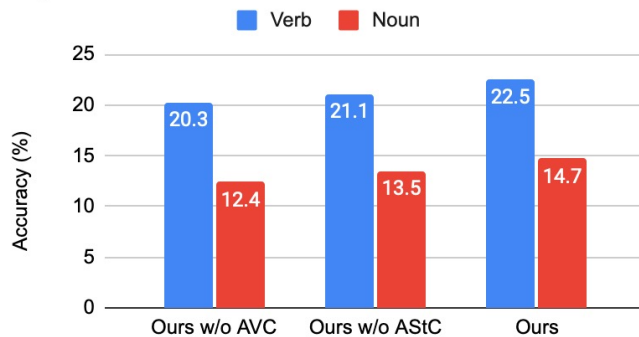
Discussion of Results (2)

Ablation Study (1) - **Our method w/o AVC and w/o AStC**: Each term enhances the representations obtained through large-scale audio-visual pre-training.

EPIC-Kitchens-100



Ego4D

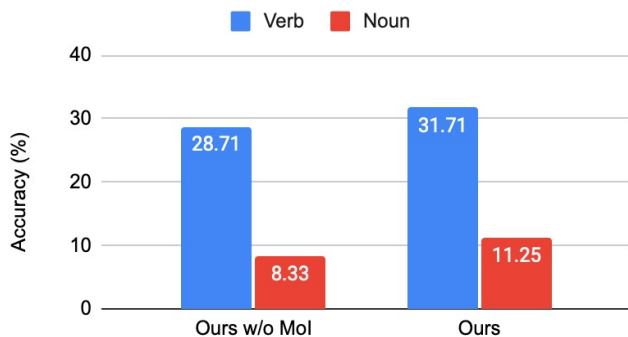


Action Recognition (Top-1 Accuracy) on EPIC-Kitchens-100 (left) and Ego4D (right)
Higher is better.

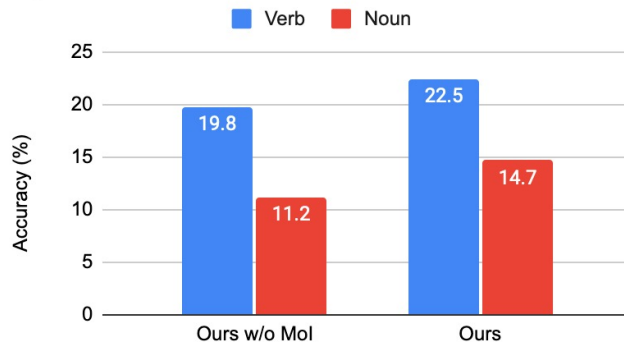
Discussion of Results (3)

Ablation Study (2) - **Our method without moments of interaction (Mol)**: Detecting moments of interaction helps representation learning.

EPIC-Kitchens-100



Ego4D



Action Recognition (Top-1 Accuracy) on EPIC-Kitchens-100 (left) and Ego4D (right)
Higher is better.

Main Takeaways/Conclusion

- We propose an audio-driven self-supervised method for learning representations of egocentric video of daily activities.
- For better representations of daily activities, learning should **focus on moments of interaction** (Mol).
 - Simple spectrogram-based Mol detector works, but there is room for improvement.
- For better representations of daily activities, models should **learn from the changes in the environment** caused by agents interacting with the world.
 - Audio can be informative of both the objects in an environment and changes in their state.