# Recognizing Manipulation Actions from State-Transformations

Nachwa ABOUBAKR, James L. CROWLEY, Rémi RONFARD

Univ. Grenoble Alpes, INRIA, Grenoble INP, CNRS, Laboratoire LIG, LJK

**CVPR** LONG BEACH CALIFORNIA June 16-20, 2019
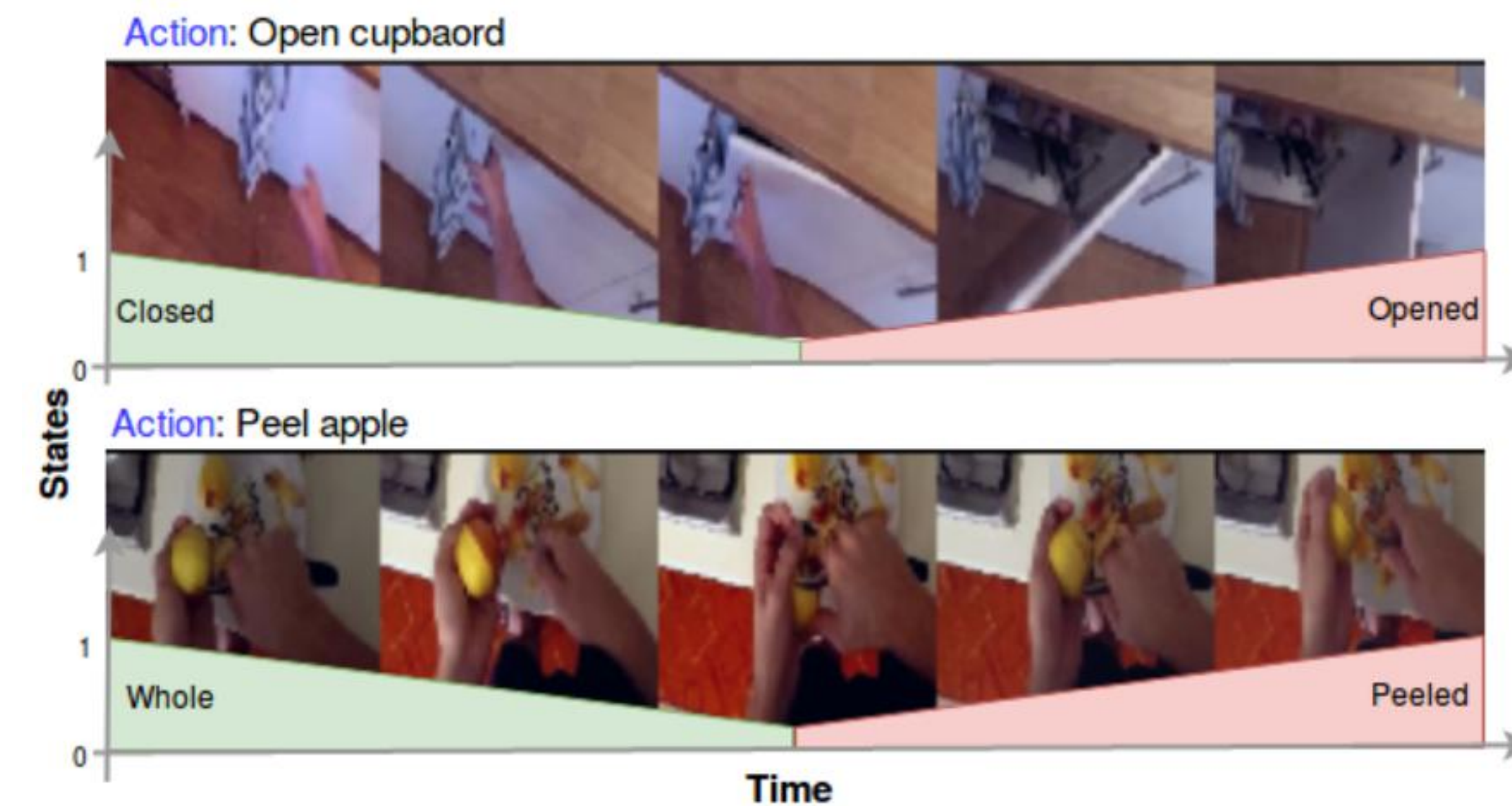
## State Transformations



Fig1: Changes in object states over time for action recognition. Two sample sequences from the EPIC kitchen dataset.

## State-Changing Actions

- State of objects are *more apparent* from still images than verbs.
- Actions can change:
  - Object's appearance,
  - Object's shape,
  - Object's position.

## State Transition

$V_i : S_{before} \rightarrow S_{after}$;   where $V_i \in$ verbs, $S_j \in$ states.
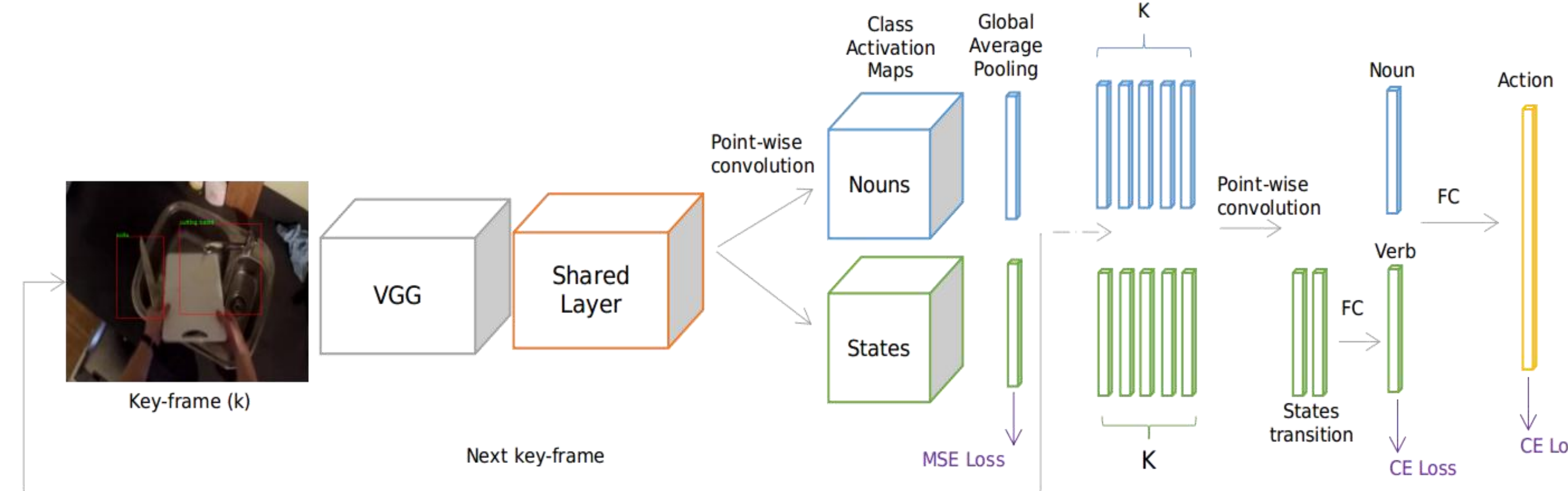
## Our Model:



Fig 2: Proposed model to learning action recognition as state transformation.

## Experiement on EPIC-Kitchens Dataset

- We defined 49 state transitions from 31 states.

### a. Action Recognition Challenge.

|  | Verbs Results | | | |
| --- | --- | --- | --- | --- |
|  | Seen kitchens subset (**S1**) | | | |
|  | Acc T1 | Acc T5 | Precision | Recall |
| Our model(RGB) | 47.41 | 81.33 | 31.20 | 20.43 |
| 2SCNN[2](RGB) | 40.44 | 83.04 | 33.74 | 15.9 |
| TSN[3](RGB) | 45.68 | 85.56 | 61.64 | 23.81 |
|  | Unseen kitchens subset (**S2**) | | | |
| Our model(RGB) | 34.35 | 69.24 | 15.09 | 11.00 |
| 2SCNN[2](RGB) | 33.12 | 73.23 | 16.06 | 9.44 |
| TSN[3](RGB) | 34.89 | 74.56 | 19.48 | 11.22 |

Table 1: Comparison of our method and baseline methods reported by [1].

### b. Results on state-changing verbs on validation set.

|  | take | put | open | close | wash | cut | mix | pour | peel | Avg |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Precision (%) | 56.7 | 59.3 | 58.8 | 39.8 | 80.1 | 74.7 | 68.9 | 39.1 | 37.7 | 57.23 |
| Recall (%) | 48.2 | 45.0 | 62.9 | 57.1 | 67.7 | 60.7 | 50.2 | 40.3 | 53.5 | 53.96 |

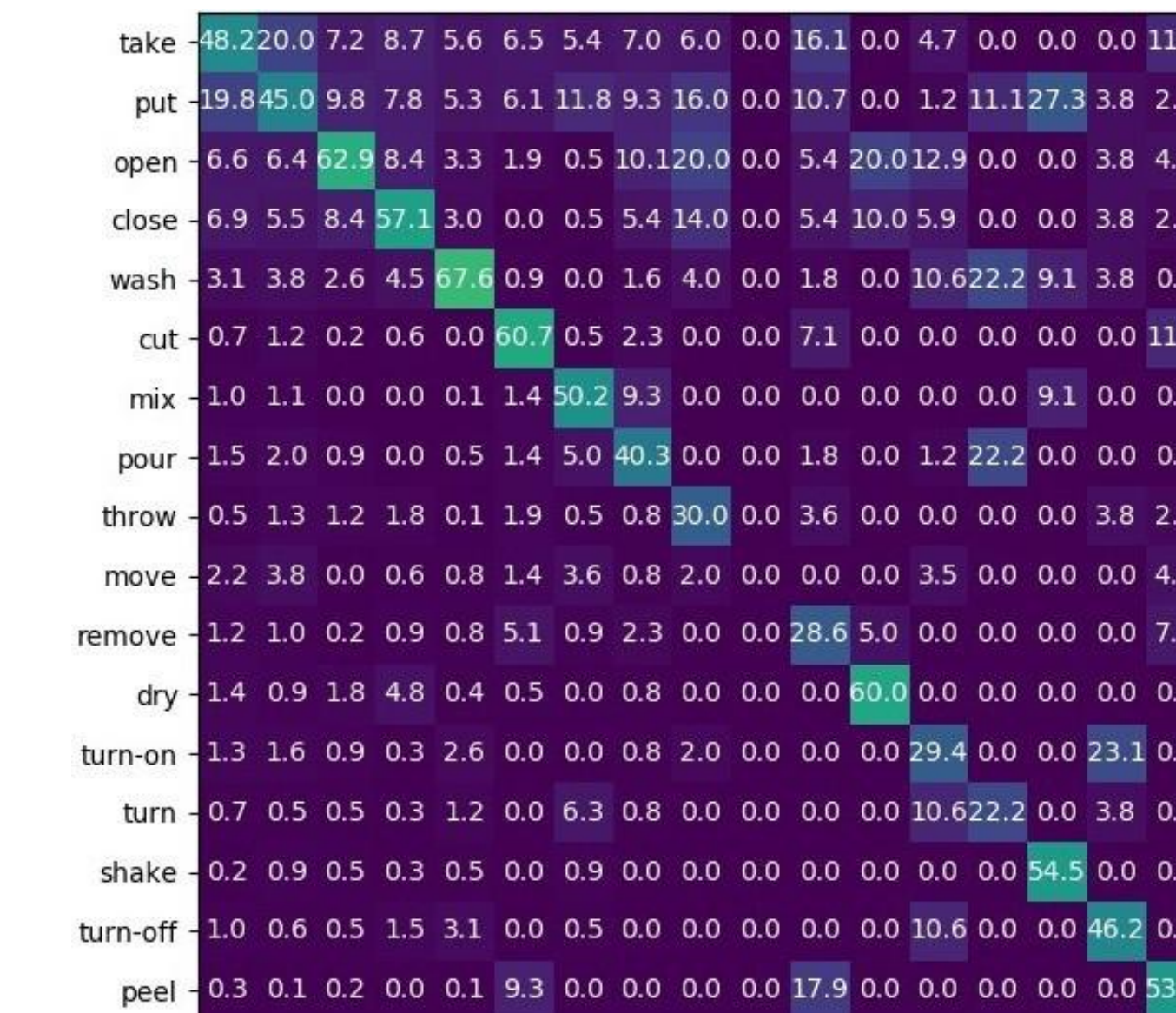Table 2: Our model performance on validation set on state-changing verbs.



Fig 3: Confusion Matrix on Validation set.

The code is available at:

## References

[1] D. Damen, H. Doughty, G. M. Farinella, S. Fidler, A. Furnari, E. Kazakos, D. Moltisanti, J. Munro, T. Perrett, W. Price, and M. Wray. Scaling egocentric vision: The epic-kitchens dataset. In European Conference on Computer Vision (ECCV), 2018.

[2] K. Simonyan and A. Zisserman. Two-stream convolutional networks for action recognition in videos. In Advances in neural information processing systems, pages 568–576, 2014.

[3] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. Van Gool. Temporal segment networks: Towards good practices for deep action recognition. In European conference on computer vision, pages 20–36. Springer, 2016.