# The Visual Centrifuge: Model-Free Layered Video Representations

## Jean-Baptiste Alayrac<sup>\*,1</sup>, Joao Carreira<sup>\*,1</sup> and Andrew Zisserman<sup>1,2</sup>

\*equal contribution, <sup>1</sup>Deepmind, <sup>2</sup>University of Oxford



Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset, Carreira and Zisserman, CVPR, 2017.

Deep multi-scale video prediction beyond mean square error, Mathieu et al, *ICLR*, 2016.

A computational approach for obstruction-free photography, Xue et al., *SIGGRAPH*, 2015.

# Experiments using the blending procedure -



### Motion vs. Static cues ablation study.



#### Qualitative samples on real world videos.





Qualitative samples on Kinetics Validation set.



Conclusion: deeper model, more outputs and predictor-corrector architecture matter.

#### Low-level vs. high-level features:

- Mixed 5c > Mixed 3c (deeper is better)
- Unmixing videos of the same class is harder (0.145 vs 0.133)
- More correlation between unmixing perf. with high level feat. than with low level feat.

#### Color experiment.







Comparison to engineered method.



#### Action recognition in challenging scenario (recognizing action in blended videos)

- I3D off-the-shelf: 22% acc.
- Centrifuge + I3D off-the-shelf: 44% acc.