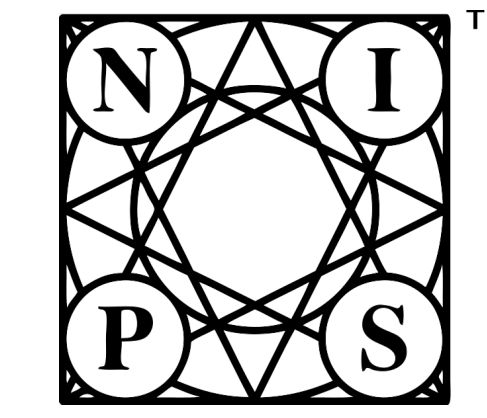




# Trajectory Convolution for Action Recognition

Yue Zhao<sup>1</sup>, Yuanjun Xiong<sup>2</sup>, Dahua Lin<sup>1</sup>

<sup>1</sup>CUHK – SenseTime Joint Lab, The Chinese University of Hong Kong <sup>2</sup>Amazon Rekognition



Neural Information  
Processing Systems  
Foundation

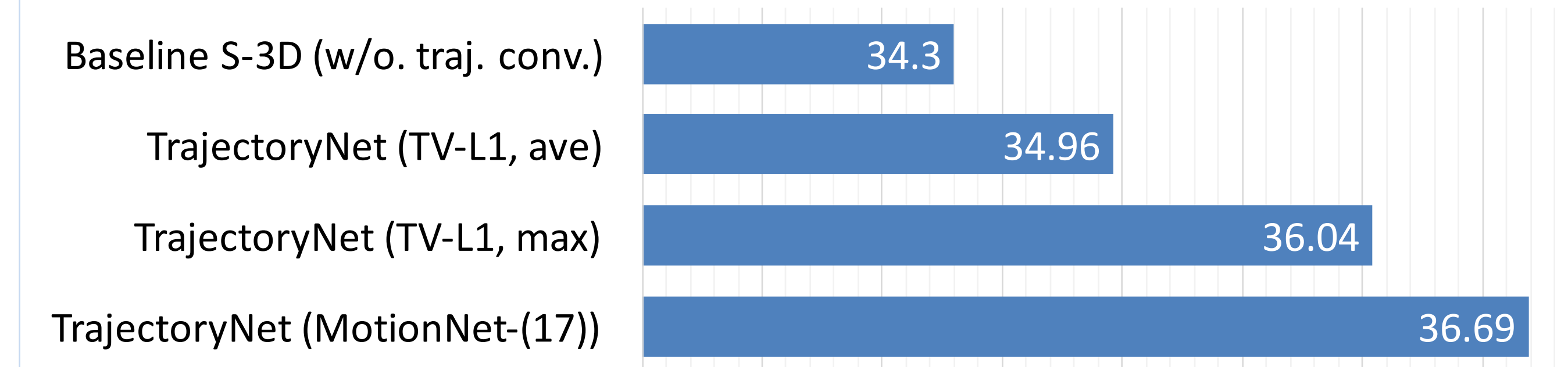
# Introduction

- How to leverage the temporal dimension is one major question in video analysis.
- The temporal convolution in Separable-3D networks, however, comes with an implicit assumption - the feature maps across time steps are well aligned.
- This assumption can be overly strong in action recognition because of *motion*.

## Comparison with previous methods

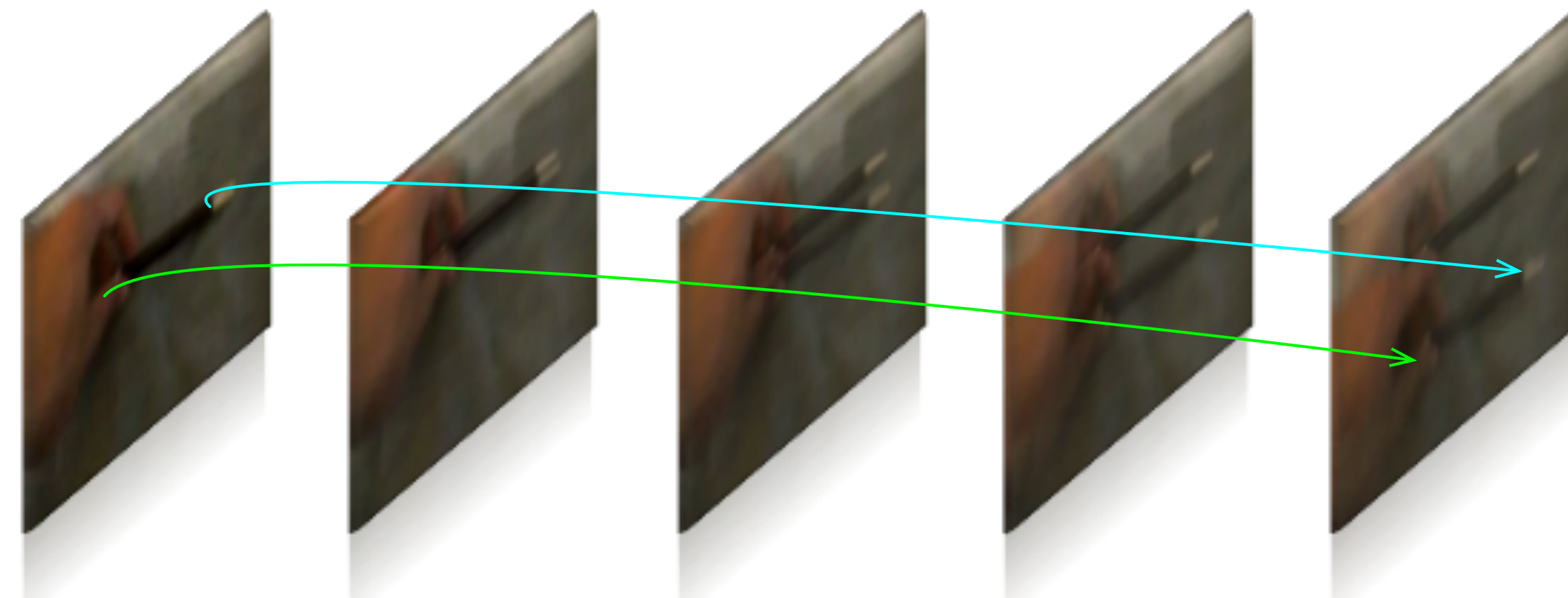
Methods	Use deep feature?	Feature tracking?	End-to-end?
STIP + HOG/HOF/MBH/...	$\times$	$\times$	$\times$
DT, iDT	$\times$	✓	$\times$
Two-stream, TSN, I3D	✓	$\times$	✓
TDD	✓	✓	$\times$
TrajectoryNet (Ours)	✓	✓	✓

## Influence of the quality of trajectory

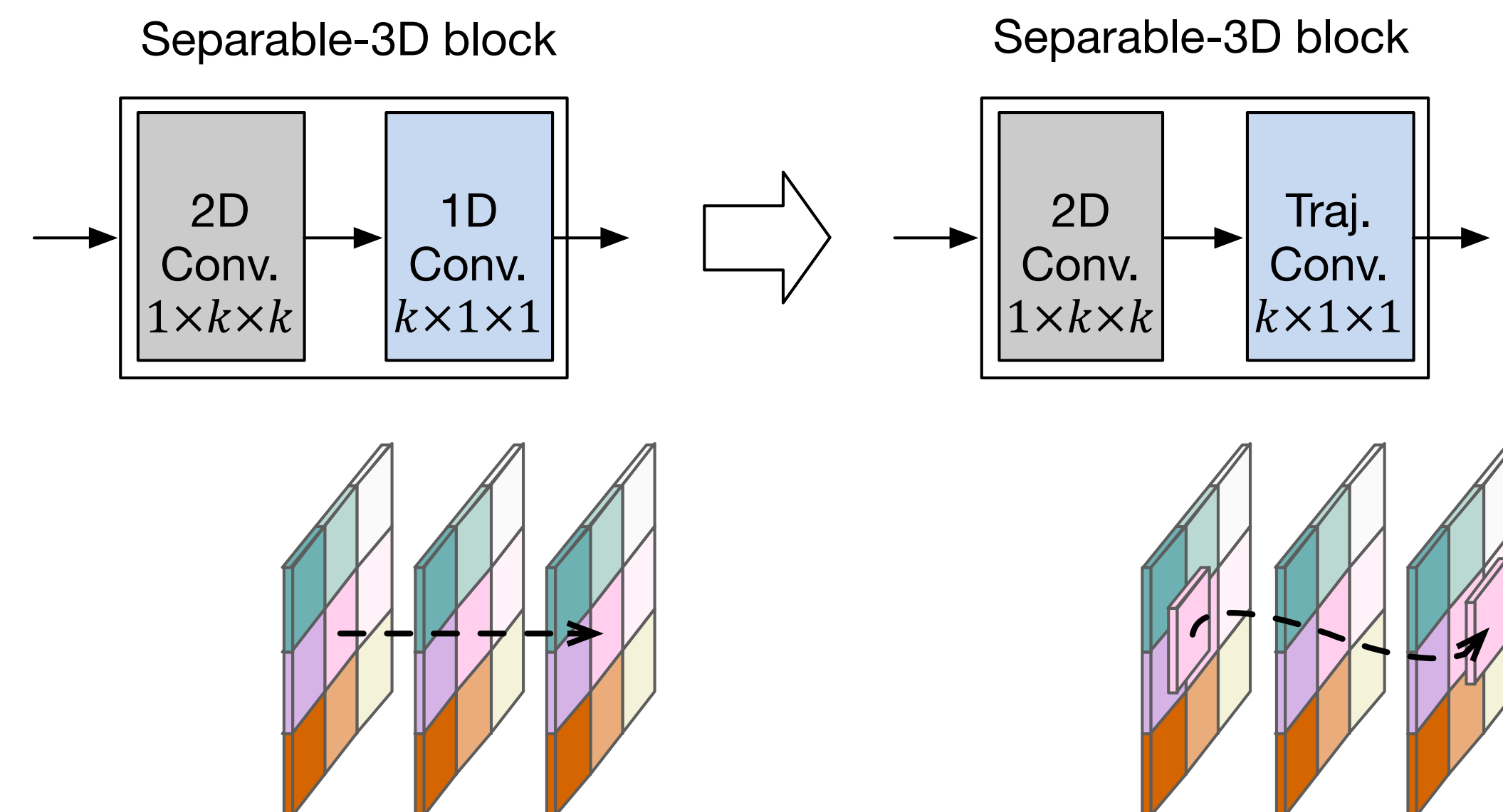


## Main idea

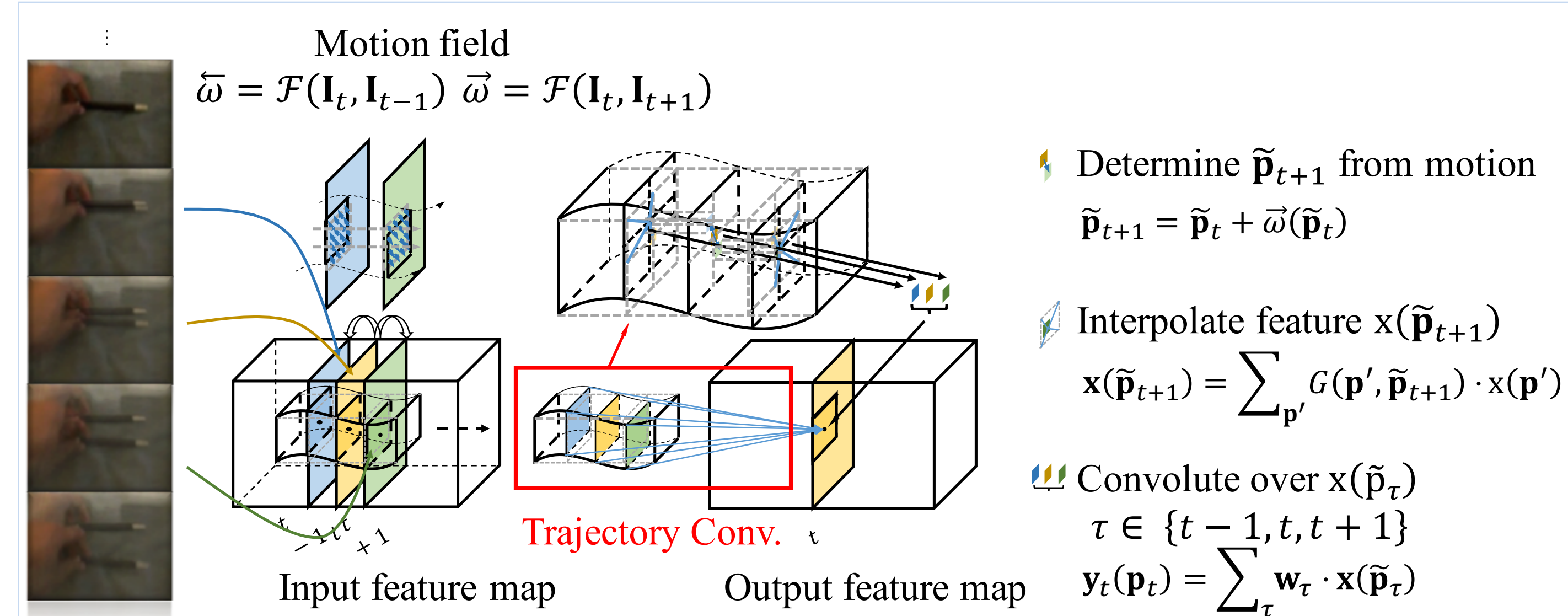
- The trajectory convolution operates along the *trajectories* that trace the pixels corresponding to the same physical points (e.g. [nib](#)), rather than at fixed pixel locations.



- The standard temporal (1D) convolution in Separable-3D can be seen as a special case of the trajectory convolution where all pixels are considered to be stationary over time.

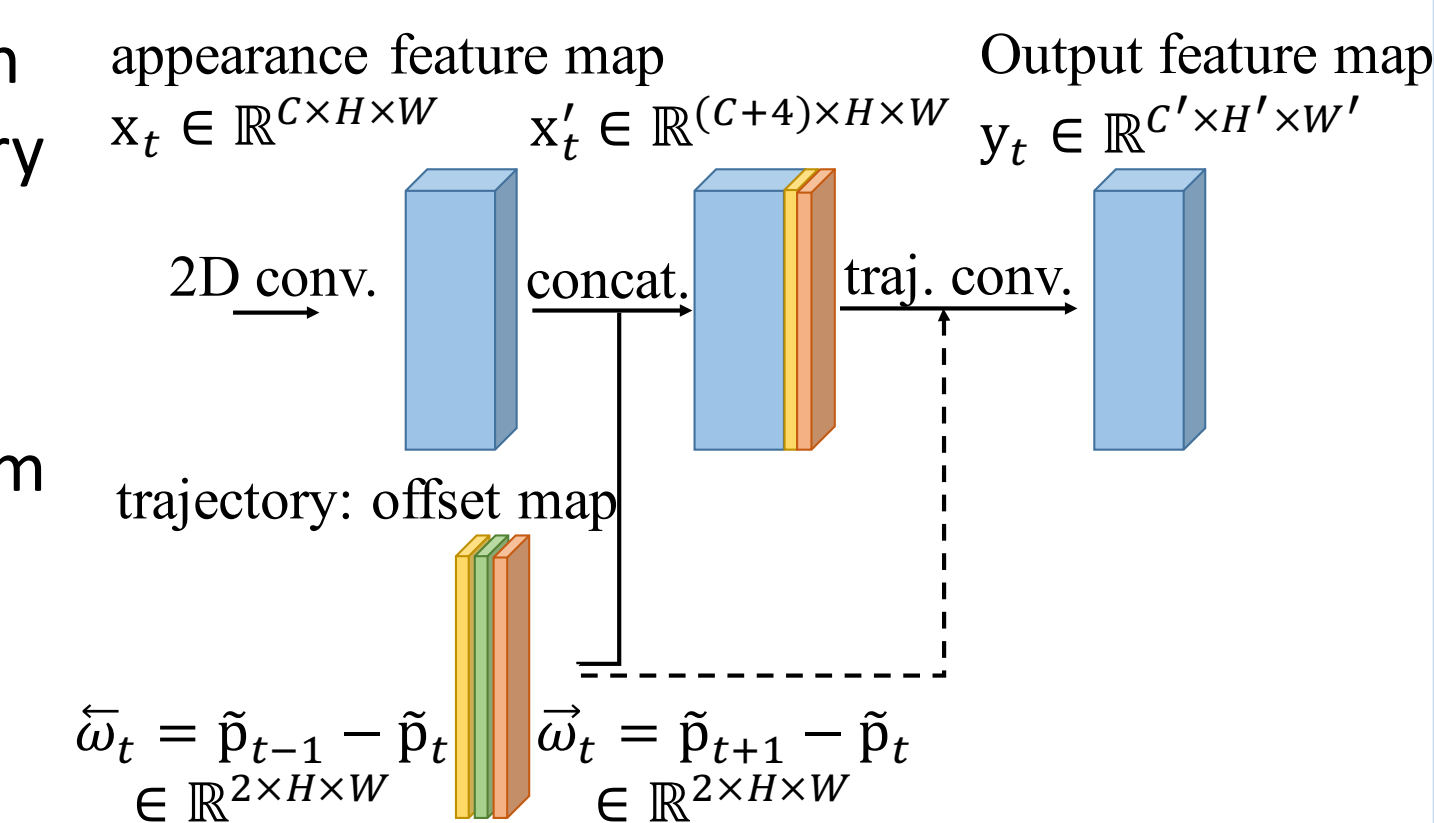


# Trajectory Convolution



# TrajectoryNet

- We describe local motion patterns at each position  $\mathbf{p}$  using the sequence of trajectory information in the form of coordinates of sampling offsets  $\{\Delta p_\tau: \tau \in [-\Delta t, \Delta t]\}$ .
- The trajectories can either be derived from pre-computed TV-L1 or be learnt from a MotionNet in an unsupervised manner.



## Main results on Something-Something-V1

Method	Backbone Network	Pre-train	Val top1
3D CNN	C3D	Sports-1M	11.5
MultiScale TRN	BN-Inception	ImageNet	34.4
ECO lite	BN-Inception + 3D-ResNet 18	Kinetics	46.4
Non-local I3D + GCN	ResNet-50	Kinetics	46.1
TrajectoryNet-MotionNet-(17) w/o. motion	ResNet-18	ImageNet	43.3
TrajectoryNet-MotionNet-(17) w/. motion	ResNet-18	ImageNet	44.0
TrajectoryNet-MotionNet-(17) w/o. motion	ResNet-18	Kinetics	47.8
TrajectoryNet-MotionNet-(17) w/. motion	ResNet-18	Kinetics	47.9

## Visualization of intermediate features

