

# 2D Deep Video Capsule Network with Temporal Shift for Action Recognition

T. Voillemin - H. Wannous - J-P. Vandeborre



# Introduction

## Motivations & challenges

- Action Recognition for Virtual/Augmented reality devices



Garcia-Hernando et al., *CVPR*, 2018

De Smedt et al., *SHREC*, 2017

- Light but efficient algorithm for real-time applications
- Working with easily accessible data as RGB and depth video stream

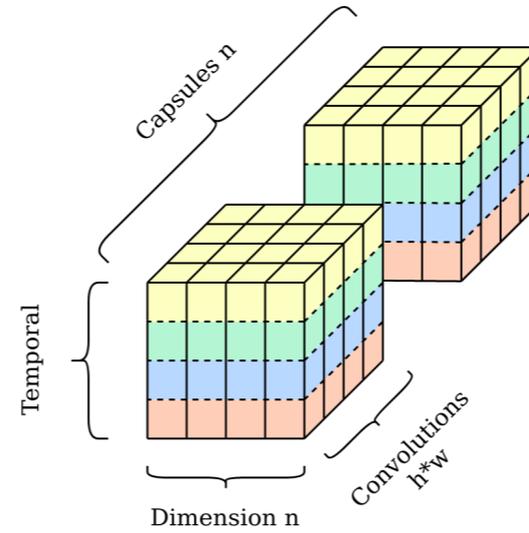
# Introduction

## Existing approaches

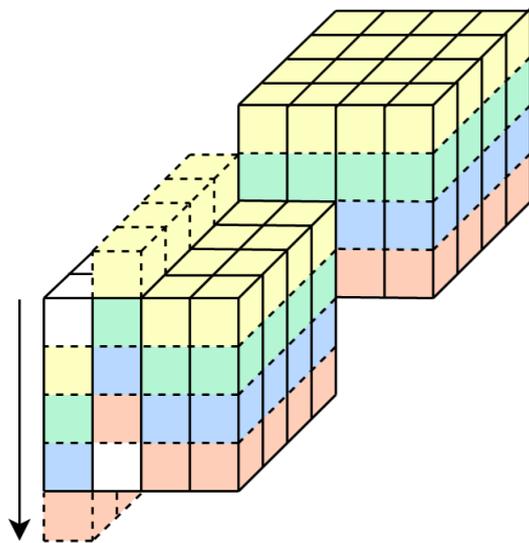
- For lightness :
  - Capsule Network (Sabour et al., *ANIPS*, 2017)
- For efficiency :
  - DeepCaps (Rajasegaran et al., *CVPR*, 2019)
  - Temporal Shift Module (Lin et al. *ICCV*, 2019)

# Proposed Method

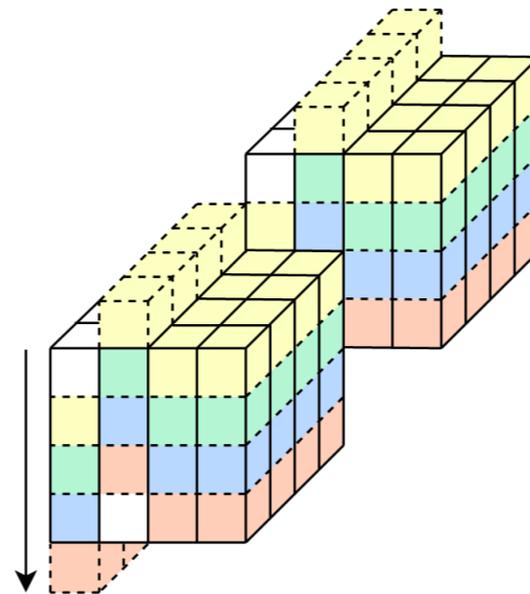
## Temporal Shift on Capsule Layer



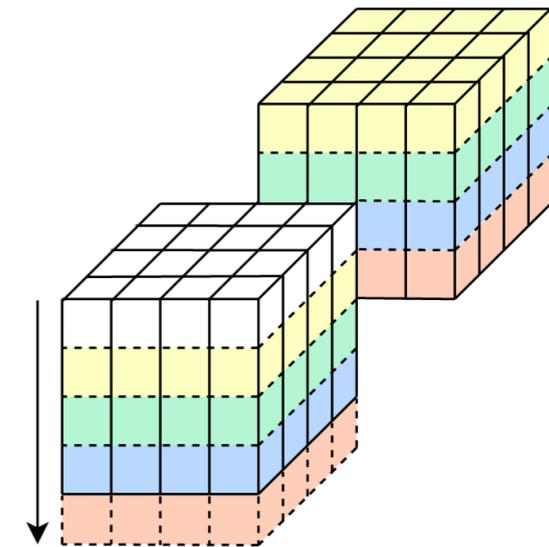
Original Capsule Layer



Shift of first kernels on first capsule



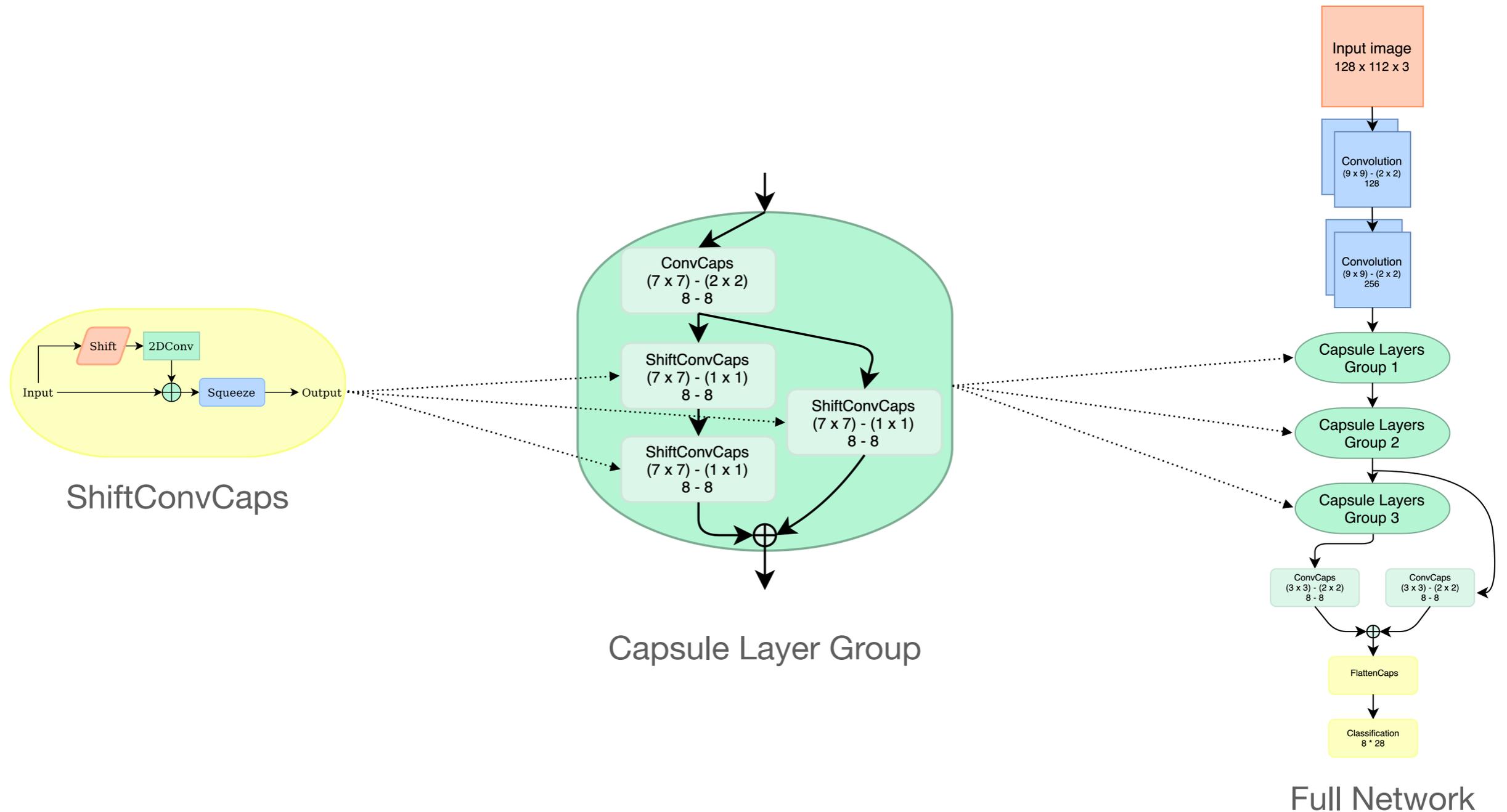
Shift of first kernels on every capsules



Shift of every kernels on first capsule

# Proposed Method

## Network Architecture



# Experiments

## Accuracy comparison over DHG dataset



De Smedt et al., *SHREC*, 2017

Method	Parameters	Accuracy (%)
TSM [2]	24 M	58.66
2D 3DCNN Fusion [3]	140 M	<b>74.41</b>
<b>Ours</b>	<b>7 M</b>	<b>68.98</b>

Comparaison on DHG28 [5]

- [1] Feichtenhofer et al., *CVPR*, 2016
- [2] Lin et al., *ICCV*, 2019
- [3] Zhang et al., *Electronics*, 2019
- [4] Garcia-Hernando et al., *CVPR*, 2018
- [5] De Smedt et al., *SHREC*, 2017

# Experiments

## Accuracy comparison over FPHA dataset

Method	Parameters	Accuracy (%)
Two stream-color [1]	46 M	61.56
Two stream-flow [1]	46 M	69.91
Two stream-all [1]	181 M	75.30
TSM [2]	24 M	71.57
<b>Ours</b>	<b>4 M</b>	<b>76.72</b>

Comparison on FPHA [4]



Garcia-Hernando et al., CVPR, 2018

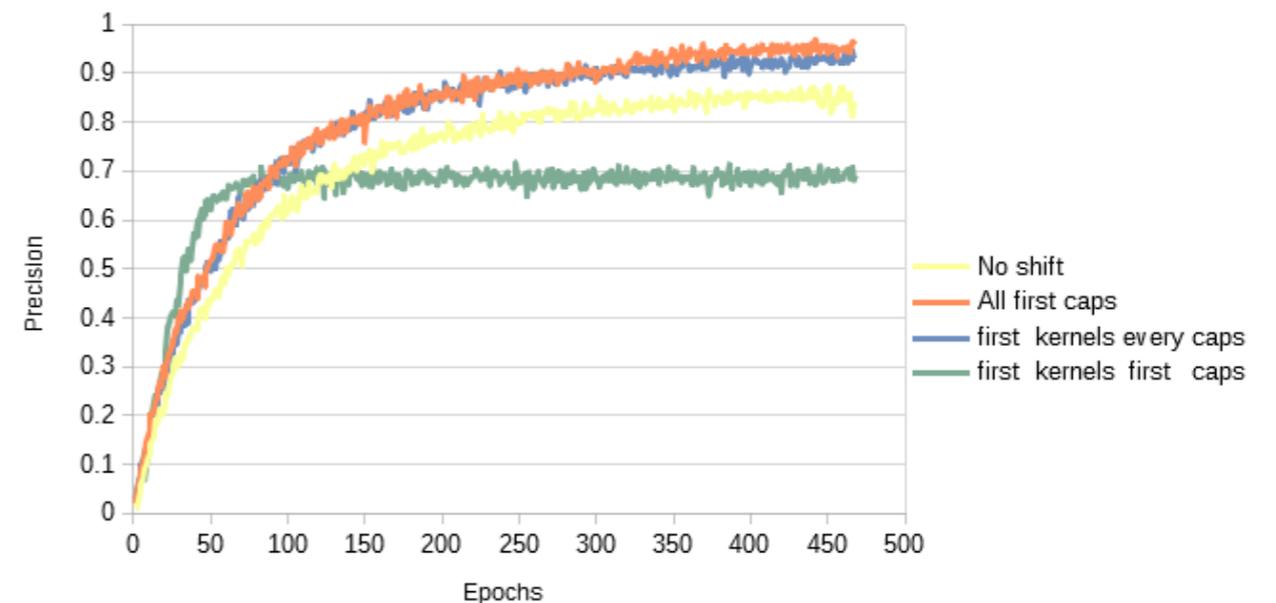
- [1] Feichtenhofer et al., CVPR, 2016
- [2] Lin et al., ICCV, 2019
- [3] Zhang et al., Electronics, 2019
- [4] Garcia-Hernando et al., CVPR, 2018
- [5] De Smedt et al., SHREC, 2017

# Experiments

## Accuracy comparison over FPHA dataset

Temporal Shift	Accuracy (%)
No shift	70.01
Shift first kernels on first capsules	64.14
Shift first kernels on every capsules	<b>74.33</b>
Shift every kernels of first capsules	<b>76.72</b>

Final testing accuracy



Training accuracy evolution over epochs

# Conclusion

## Contributions & perspectives

- First 2D Capsule Network for video understanding
- Implementation of temporal shift over capsule layer
- Outperform or near state-of-the-art with 10 to 40 times less parameters
  
- Adapt capsules for hand skeleton data
- Use it in real-time directly on augmented/virtual reality devices

# Thank you for your attention

T. Voillemin - H. Wannous - J-P. Vandeborre

