

# Gate-Shift Networks for Video Action Recognition

Swathikiran Sudhakaran<sup>1</sup>

Sergio Escalera<sup>2,3</sup>

Oswald Lanz<sup>1</sup>

<sup>1</sup>Fondazione Bruno Kessler, Italy

<sup>2</sup>Computer Vision Center, Spain

<sup>3</sup>Universitat de Barcelona, Spain



# Motivation

Video action recognition requires spatio-temporal reasoning

Putting something similar to other things that are already on the table



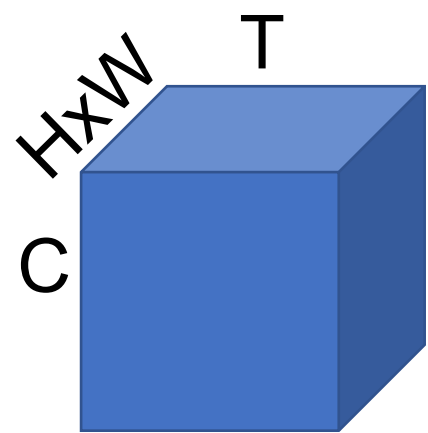
Taking one of many similar things on the table

# Contribution

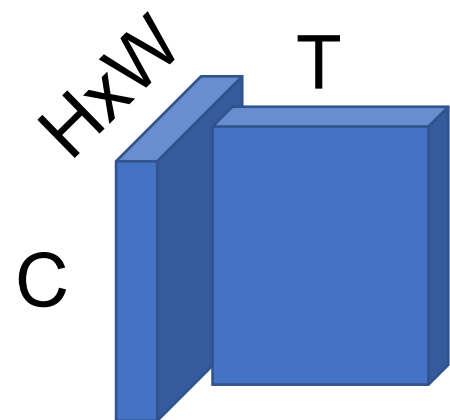
Large number of parameters in 3D CNNs require **large scale annotated data** for training

Existing approaches address this problem by a **hard-wired** decomposition of the 3D kernels which is **suboptimal**

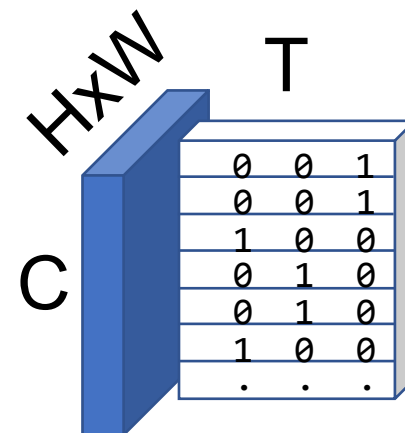
**GSM** leverages spatial gating for **adaptive feature propagation**



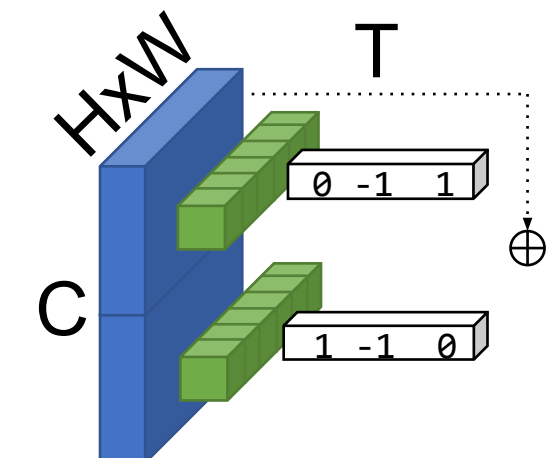
C3D



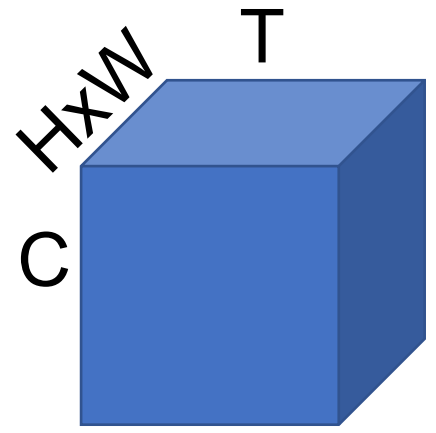
S3D /  
R(2+1)D



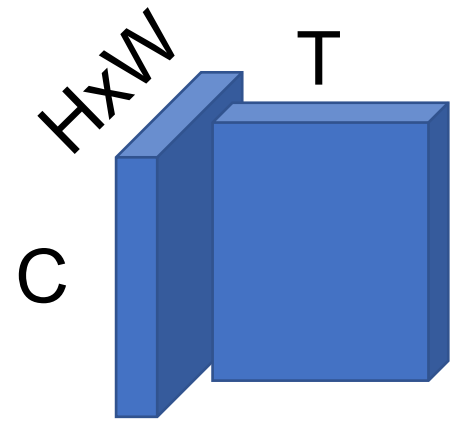
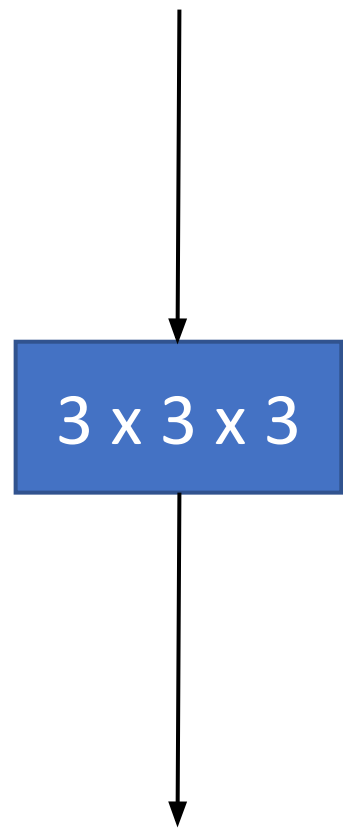
TSM



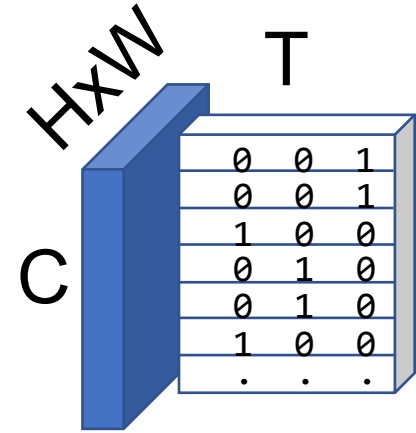
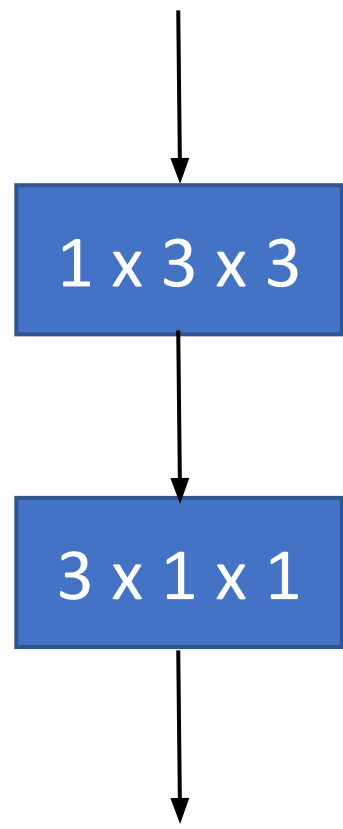
**GSM**



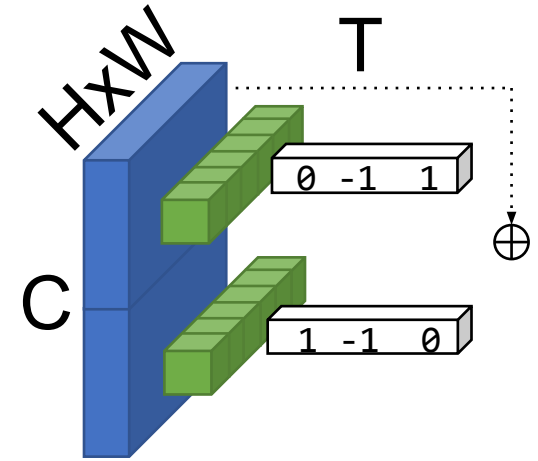
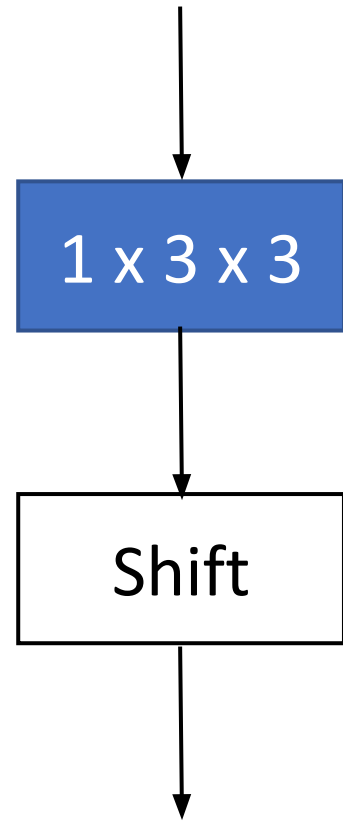
C3D



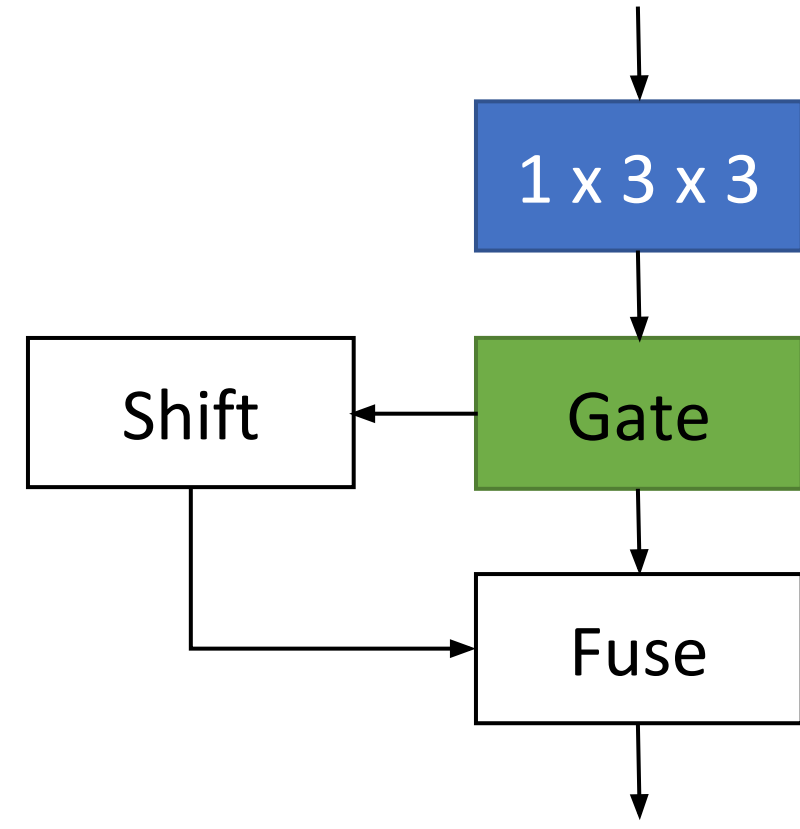
S3D /  
R(2+1)D



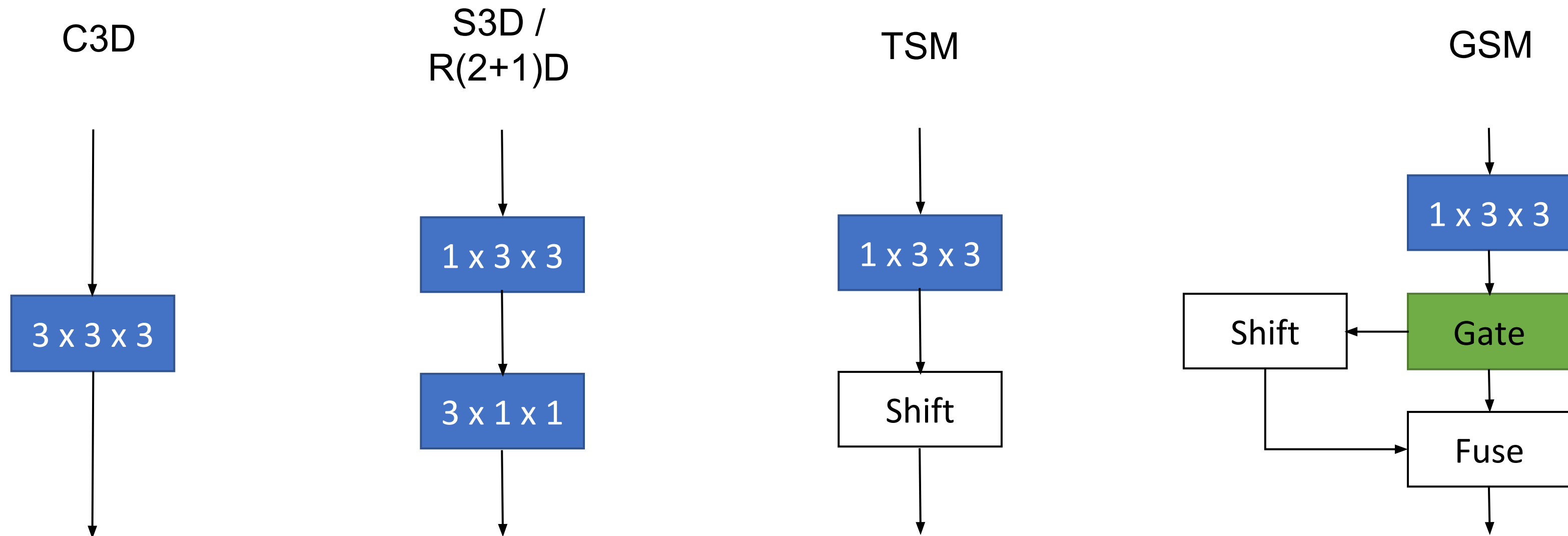
TSM



GSM

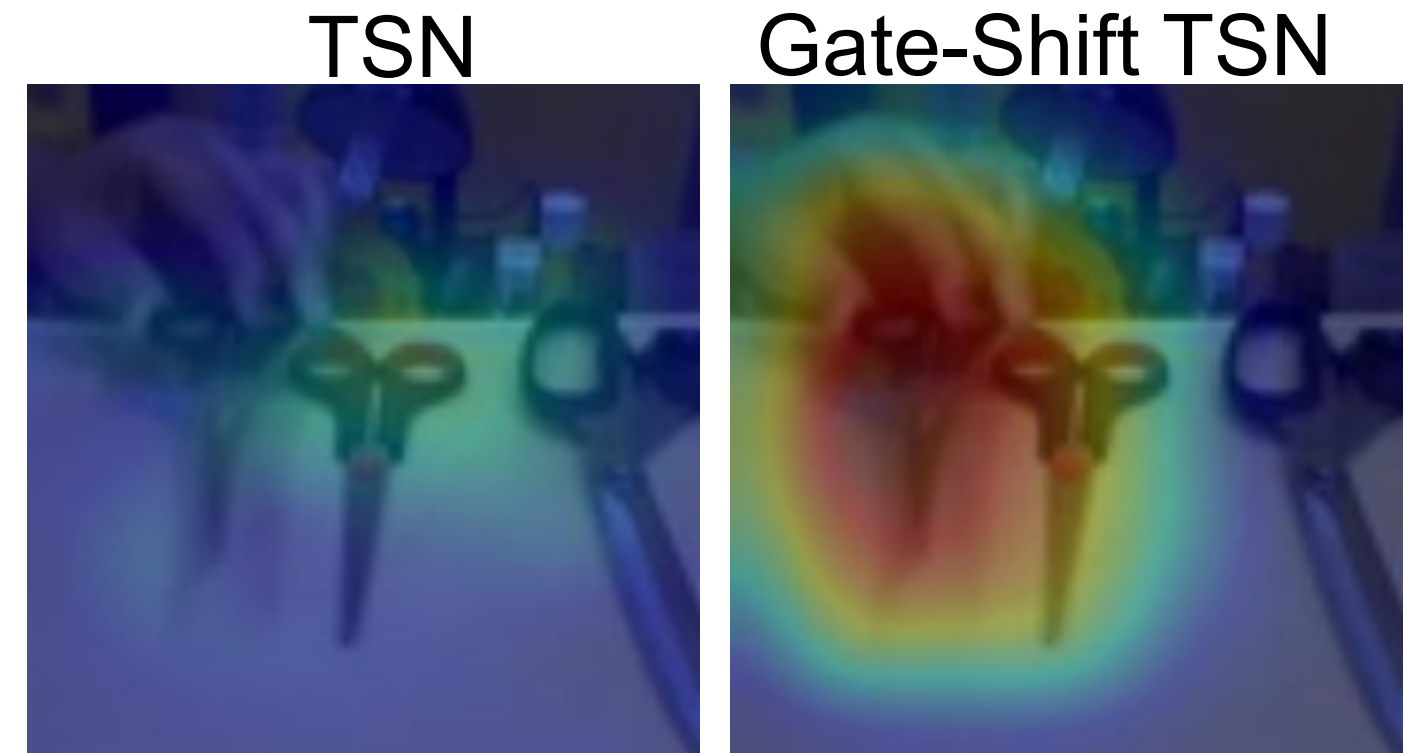
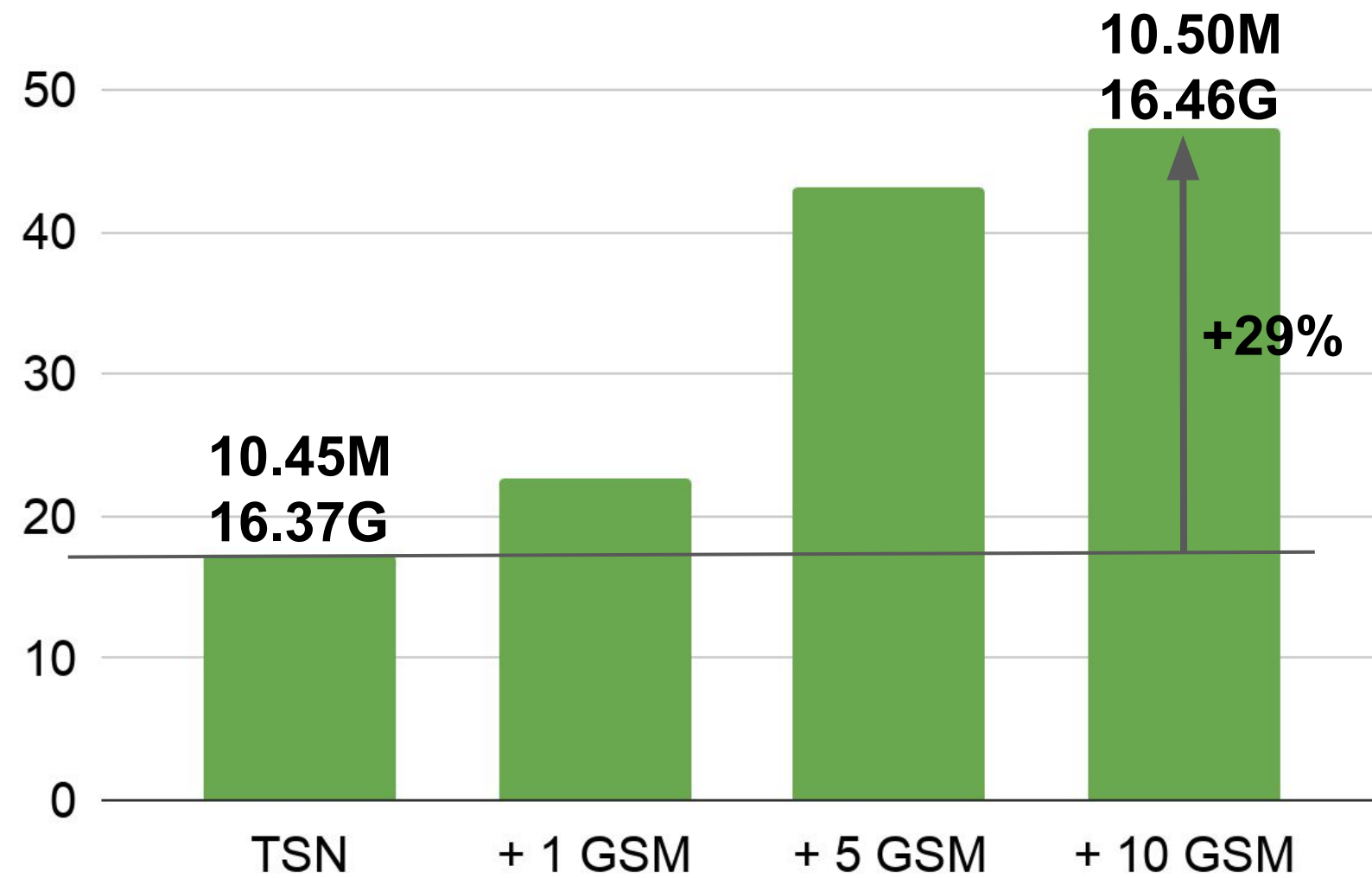


GSM develops a **flexible** and **data dependent decomposition** of 3D kernels with **reduced parameters** and **computational overhead**

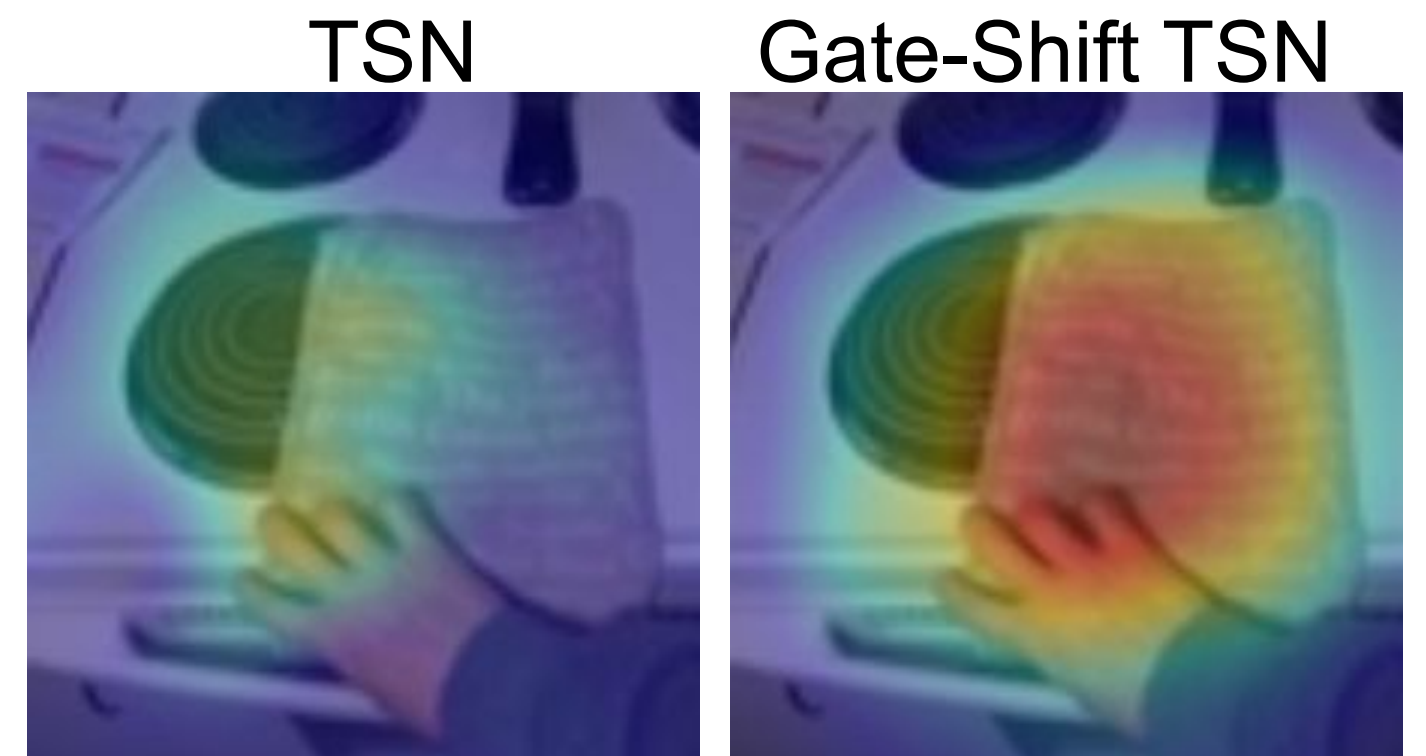


# Effectiveness of GSM

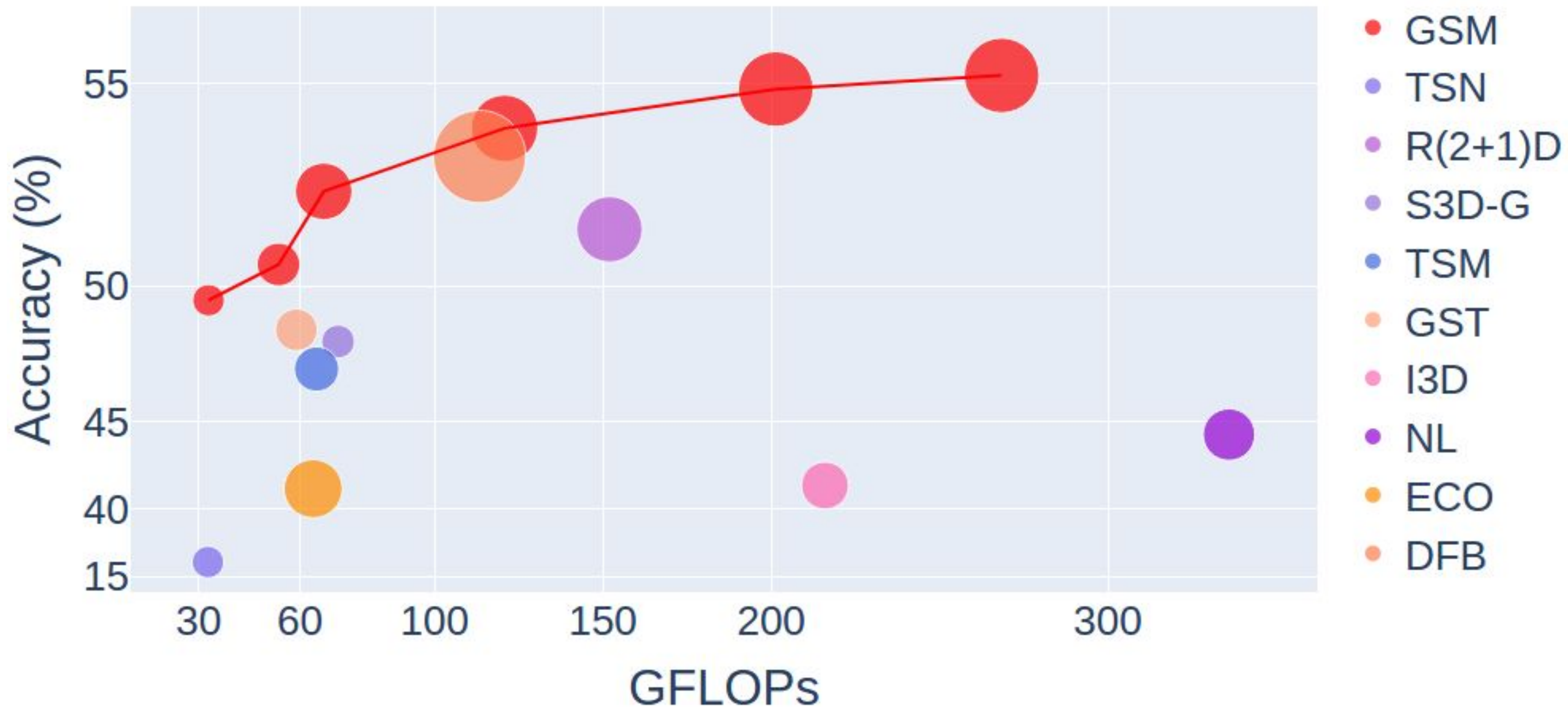
Ablation study on Sth-V1



Putting sth similar to other things that are already on the table



Unfolding sth



State-of-the-art recognition accuracy of **55%** on Something Something-V1

# Gate-Shift Networks for Video Action Recognition

Swathikiran Sudhakaran<sup>1</sup>

Sergio Escalera<sup>2,3</sup>

Oswald Lanz<sup>1</sup>

<sup>1</sup>Fondazione Bruno Kessler, Italy

<sup>2</sup>Computer Vision Center, Spain

<sup>3</sup>Universitat de Barcelona, Spain

