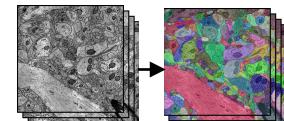


Parallel Separable 3D Convolution for Video and Volumetric Data Understanding

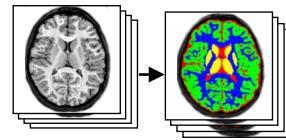
Felix Gonda, Donglai Wei, Toufiq Parag, Hanspeter Pfister



Action Rec. in Videos



Neuron Segmentation



MRI Segmentation

Motivation

3D convolutions capture spatial-temporal signal in video and volumetric data.

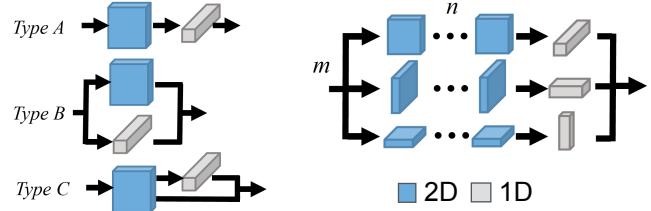
But are expensive:

- Computation speed
- Memory constraints

We use tensor factorization to improve speed and reduce model size.

Our Approach

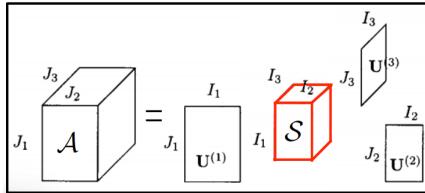
Replace 3D convolution with computational graph of orthogonal pairs of 2D and 1D



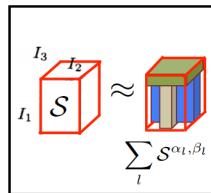
Method	# Streams	# 2D Conv.	Type
P3D ⁽⁺⁾	1	1	ABC
R(2+1)D ^(#)	1	1	A
P _m SC _n (ours)	m	n	A

(+) Pseudo 3D (#) Facebook

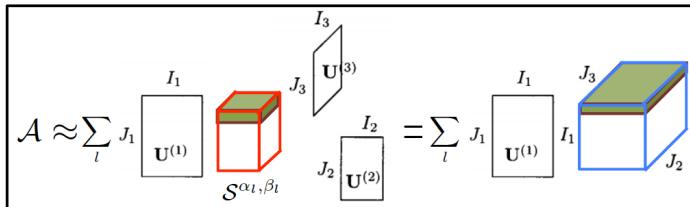
Parallel Streams



Decompose A into orthogonal matrices U(k)

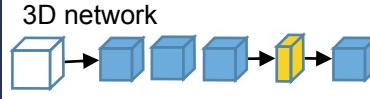


Sparse approximation

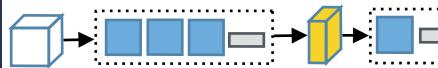


Re-group into (2+1)D (different 2D non-zero planes)

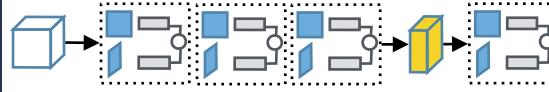
Approximation Examples



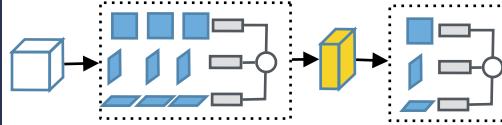
(a) Single stream group replacement



(b) Dual stream 1:1 replacement



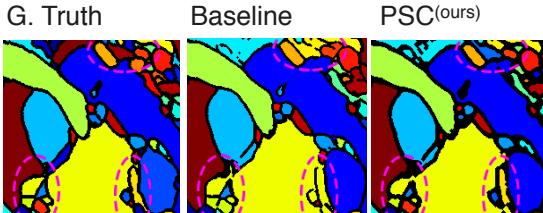
(c) Triple stream group replacement



Legend:
█ Data
█ 3D Conv.
█ Pooling
○ Concat

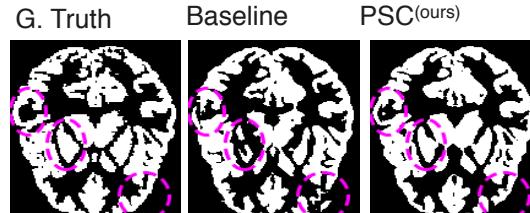
Results

Neuron Segmentation
(U-Net on FIBSEM)



G. Truth Baseline PSC^(ours)

MRI Segmentation
(DenseNet on IBSR)



Action Recognition
(ResNet on UCF101)

Type	VI	# Param
3D	0.10, 0.48	21M
P ₁ SC ₁ ^(ours)	0.11, 0.24	11M
P ₂ SC ₂ ^(ours)	0.07, 0.23	12M
P ₂ SC ₃ ^(ours)	0.08, 0.28	10M

Type	WM	GM	CSF	#Param
3D	91.3	91.6	94.7	1.6M
P ₁ SC ₁ ^(#)	95.1	94.1	93.2	4.7M
P ₁ SC ₁ ^(ours)	95.7	96.1	96.3	2.5M
P ₂ SC ₂ ^(ours)	95.2	96.1	97.6	1.4M

Type	Acc.	#Param
P3D ⁽⁺⁾	88.6	261M
3D	85.4	64M
P ₁ SC ₁ ^(#)	93.6	39M
P ₁ SC ₁ ^(ours)	89.7	39M
P ₂ SC ₂ ^(ours)	92.3	49M
P ₂ SC ₃ ^(ours)	91.4	39M